



August 1998

A Lexicalized Tree Adjoining Grammar for English

The XTAG Research Group
University of Pennsylvania

Follow this and additional works at: http://repository.upenn.edu/ircs_reports

Research Group, The XTAG, "A Lexicalized Tree Adjoining Grammar for English" (1998). *IRCS Technical Reports Series*. 63.
http://repository.upenn.edu/ircs_reports/63

University of Pennsylvania Institute for Research in Cognitive Science Technical Report No. IRCS-98-18.

This paper is posted at ScholarlyCommons. http://repository.upenn.edu/ircs_reports/63
For more information, please contact libraryrepository@pobox.upenn.edu.

A Lexicalized Tree Adjoining Grammar for English

Abstract

This document describes a sizable grammar of English written in the TAG formalism and implemented for use with the XTAG system. This report and the grammar described herein supersedes the TAG grammar described in [XTAG-Group,1995]. The English grammar described in this report is based on the TAG formalism developed in [Joshi *et al.*, 1975], which has been extended to include lexicalization ([Schabes *et al.*, 1998]) and unification-based feature structures ([Vijay-Shanker and Joshi, 1991]). The range of syntactic phenomena that can be handled is large and includes auxiliaries (including inversion), copula, raising and small clause constructions, topicalization, relative clauses, infinitives, gerunds, passives, adjuncts, it-clefts, wh-clefts, PRO constructions, noun-noun modifications, extraposition, determiner sequences, genitives, negation, noun-verb contractions, sentential adjuncts and imperatives. This technical report corresponds to the XTAG Release 8/31/98. The XTAG grammar is continuously updated with the addition of new analyses and modification of old ones, and an online version of this report can be found at the XTAG web page: <http://www.cis.upenn.edu/~xtag/>.

Comments

University of Pennsylvania Institute for Research in Cognitive Science Technical Report No. IRCS-98-18.

A Lexicalized Tree Adjoining Grammar for English

The XTAG Research Group

Institute for Research in Cognitive Science
University of Pennsylvania
3401 Walnut St., Suite 400A
Philadelphia, PA 19104-6228

<http://www.cis.upenn.edu/~xtag>

August 31, 1998

Contents

I	General Information	1
1	Getting Around	3
2	Feature-Based, Lexicalized Tree Adjoining Grammars	5
2.1	TAG formalism	5
2.2	Lexicalization	7
2.3	Unification-based features	8
3	Overview of the XTAG System	11
3.1	System Description	11
3.1.1	Tree Selection	11
3.1.2	Tree Database	13
3.1.3	Tree Grafting	14
3.1.4	The Grammar Development Environment	14
3.2	Computer Platform	17
4	Underview	19
4.1	Subcategorization Frames	19
4.2	Complements and Adjuncts	20
4.3	Non-S constituents	21
4.4	Case Assignment	21
4.4.1	Approaches to Case	21
4.4.1.1	Case in GB theory	21
4.4.1.2	Minimalism and Case	22
4.4.2	Case in XTAG	22
4.4.3	Case Assigners	23
4.4.3.1	Prepositions	23
4.4.3.2	Verbs	23
4.4.4	PRO in a unification based framework	25
II	Verb Classes	29
5	Where to Find What	31

6	Verb Classes	38
6.1	Intransitive: Tnx0V	38
6.2	Transitive: Tnx0Vnx1	39
6.3	Ditransitive: Tnx0Vnx1nx2	40
6.4	Ditransitive with PP: Tnx0Vnx1pnx2	41
6.5	Ditransitive with PP shift: Tnx0Vnx1tonx2	42
6.6	Sentential Complement with NP: Tnx0Vnx1s2	43
6.7	Intransitive Verb Particle: Tnx0Vpl	44
6.8	Transitive Verb Particle: Tnx0Vplnx1	45
6.9	Ditransitive Verb Particle: Tnx0Vplnx1nx2	46
6.10	Intransitive with PP: Tnx0Vpnx1	47
6.11	Predicative Multi-word with Verb, Prep anchors: Tnx0VPnx1	48
6.12	Sentential Complement: Tnx0Vs1	49
6.13	Intransitive with Adjective: Tnx0Vax1	49
6.14	Transitive Sentential Subject: Ts0Vnx1	50
6.15	Light Verbs: Tnx0IVN1	51
6.16	Ditransitive Light Verbs with PP Shift: Tnx0IVN1Pnx2	51
6.17	NP It-Cleft: TItVnx1s2	53
6.18	PP It-Cleft: TItVpnx1s2	53
6.19	Adverb It-Cleft: TItVad1s2	54
6.20	Adjective Small Clause Tree: Tnx0Ax1	54
6.21	Adjective Small Clause with Sentential Complement: Tnx0A1s1	55
6.22	Adjective Small Clause with Sentential Subject: Ts0Ax1	56
6.23	Equative <i>BE</i> : Tnx0BEnx1	57
6.24	NP Small Clause: Tnx0N1	57
6.25	NP Small Clause with Sentential Complement: Tnx0N1s1	58
6.26	NP Small Clause with Sentential Subject: Ts0N1	58
6.27	PP Small Clause: Tnx0Pnx1	59
6.28	Exhaustive PP Small Clause: Tnx0Px1	60
6.29	PP Small Clause with Sentential Subject: Ts0Pnx1	60
6.30	Intransitive Sentential Subject: Ts0V	61
6.31	Sentential Subject with 'to' complement: Ts0Vtonx1	62
6.32	PP Small Clause, with Adv and Prep anchors: Tnx0ARBPnx1	62
6.33	PP Small Clause, with Adj and Prep anchors: Tnx0APnx1	63
6.34	PP Small Clause, with Noun and Prep anchors: Tnx0NPnx1	64
6.35	PP Small Clause, with Prep anchors: Tnx0PPnx1	65
6.36	PP Small Clause, with Prep and Noun anchors: Tnx0PNaPnx1	65
6.37	PP Small Clause with Sentential Subject, and Adv and Prep anchors: Ts0ARBPnx1	66
6.38	PP Small Clause with Sentential Subject, and Adj and Prep anchors: Ts0APnx1	67
6.39	PP Small Clause with Sentential Subject, and Noun and Prep anchors: Ts0NPnx1	67
6.40	PP Small Clause with Sentential Subject, and Prep anchors: Ts0PPnx1	69
6.41	PP Small Clause with Sentential Subject, and Prep and Noun anchors: Ts0PNaPnx1	69
6.42	Predicative Adjective with Sentential Subject and Complement: Ts0A1s1	70
6.43	Locative Small Clause with Ad anchor: Tnx0nx1ARB	70
6.44	Exceptional Case Marking: TXnx0Vs1	72

6.45	Idiom with V, D, and N anchors: Tnx0VDN1	73
6.46	Idiom with V, D, A, and N anchors: Tnx0VDAN1	73
6.47	Idiom with V and N anchors: Tnx0VN1	74
6.48	Idiom with V, A, and N anchors: Tnx0VAN1	75
6.49	Idiom with V, D, A, N, and Prep anchors: Tnx0VDAN1Pnx2	76
6.50	Idiom with V, A, N, and Prep anchors: Tnx0VAN1Pnx2	76
6.51	Idiom with V, N, and Prep anchors: Tnx0VN1Pnx2	77
6.52	Idiom with V, D, N, and Prep anchors: Tnx0VDN1Pnx2	78
7	Ergatives	79
8	Sentential Subjects and Sentential Complements	81
8.1	S or VP complements?	81
8.2	Complementizers and Embedded Clauses in English: The Data	82
8.3	Features Required	84
8.4	Distribution of Complementizers	84
8.5	Case assignment, <i>for</i> and the two <i>to</i> 's	85
8.6	Sentential Complements of Verbs	86
8.6.1	Exceptional Case Marking Verbs	88
8.7	Sentential Subjects	91
8.8	Nouns and Prepositions taking Sentential Complements	92
8.9	PRO control	93
8.9.1	Types of control	93
8.9.2	A feature-based analysis of PRO control	93
8.9.3	The nature of the control feature	94
8.9.4	Long-distance transmission of control features	94
8.9.5	Locality constraints on control	95
8.10	Reported speech	96
9	The English Copula, Raising Verbs, and Small Clauses	97
9.1	Usages of the copula, raising verbs, and small clauses	97
9.1.1	Copula	97
9.1.2	Raising Verbs	98
9.1.3	Small Clauses	99
9.1.4	Raising Adjectives	99
9.2	Various Analyses	99
9.2.1	Main Verb Raising to INFL + Small Clause	99
9.2.2	Auxiliary + Null Copula	100
9.2.3	Auxiliary + Predicative Phrase	100
9.2.4	Auxiliary + Small Clause	100
9.3	XTAG analysis	101
9.4	Non-predicative <i>BE</i>	104
10	Ditransitive constructions and dative shift	107
11	It-clefts	110

III Sentence Types	113
12 Passives	115
13 Extraction	117
13.1 Topicalization and the value of the <inv> feature	119
13.2 Extracted subjects	119
13.3 Wh-moved NP complement	120
13.4 Wh-moved object of a P	121
13.5 Wh-moved PP	123
13.6 Wh-moved S complement	123
13.7 Wh-moved Adjective complement	124
14 Relative Clauses	125
14.1 Complementizers and clauses	127
14.1.1 Further constraints on the null Comp ϵ_C	130
14.2 Reduced Relatives	130
14.2.1 Restrictive vs. Non-restrictive relatives	131
14.3 External syntax	131
14.4 Other Issues	132
14.4.1 Interaction with adjoined Comps	132
14.4.2 Adjunction on PRO	132
14.4.3 Adjunct relative clauses	133
14.4.4 ECM	133
14.5 Cases not handled	133
14.5.1 Partial treatment of free-relatives	133
14.5.2 Adjunct P-stranding	134
14.5.3 Overgeneration	134
14.5.3.1 <i>how</i> as <i>wh</i> -NP	134
14.5.3.2 <i>for</i> -trace effects	134
14.5.3.3 Internal head constraint	135
14.5.3.4 Overt Comp constraint on stacked relatives	135
15 Adjunct Clauses	136
15.0.4 Multi-word Subordinating Conjunctions	137
15.1 “Bare” Adjunct Clauses	137
15.2 Discourse Conjunction	138
16 Imperatives	140
17 Gerund NP’s	142
17.1 Determiner Gerunds	143
17.2 NP Gerunds	144
17.3 Gerund Passives	146

IV Other Constructions	149
18 Determiners and Noun Phrases	151
18.1 The Wh-Feature	156
18.2 Multi-word Determiners	156
18.3 Genitive Constructions	158
18.4 Partitive Constructions	159
18.5 Adverbs, Noun Phrases, and Determiners	161
19 Modifiers	165
19.1 Adjectives	165
19.2 Noun-Noun Modifiers	167
19.3 Time Noun Phrases	169
19.4 Prepositions	170
19.5 Adverbs	173
19.6 Locative Adverbial Phrases	179
20 Auxiliaries	183
20.1 Non-inverted sentences	184
20.2 Inverted Sentences	187
20.3 Do-Support	189
20.3.1 In negated sentences	190
20.3.2 In inverted yes/no questions	191
20.4 Infinitives	191
20.5 Semi-Auxiliaries	192
20.5.1 Marginal Modal <i>dare</i>	192
20.5.2 Other semi-auxiliaries	192
20.5.3 Other Issues	193
21 Conjunction	194
21.1 Introduction	194
21.2 Adjective, Adverb, Preposition and PP Conjunction	194
21.3 Noun Phrase and Noun Conjunction	194
21.4 Determiner Conjunction	195
21.5 Sentential Conjunction	195
21.6 Comma as a conjunction	196
21.7 <i>But-not</i> , <i>not-but</i> , <i>and-not</i> and ϵ - <i>not</i>	199
21.8 <i>To</i> as a Conjunction	200
21.9 Predicative Coordination	200
21.10 Pseudo-coordination	204
22 Comparatives	205
22.1 Introduction	205
22.2 Metalinguistic Comparatives	205
22.3 Propositional Comparatives	208
22.3.1 Nominal Comparatives	208

22.3.2	Adjectival Comparatives	213
22.3.3	Adverbial Comparatives	216
22.4	Future Work	218
23	Punctuation Marks	219
23.1	Appositives, parentheticals and vocatives	220
23.1.1	$\beta_{nx}PU_{nx}PU$	220
23.1.2	$\beta_nPU_{nx}PU$	221
23.1.3	$\beta_{nx}PU_{nx}$	221
23.1.4	$\beta_{PU_{px}PU_{vx}}$	222
23.1.5	$\beta_{puARB_{pux}}$	222
23.1.6	β_sPU_{nx}	224
23.1.7	$\beta_{nx}PU_s$	225
23.2	Bracketing punctuation	225
23.2.1	Simple bracketing	225
23.2.2	β_sPU_sPU	227
23.3	Punctuation trees containing no lexical material	227
23.3.1	αPU	227
23.3.2	βPU_s	227
23.3.3	β_sPU_s	229
23.3.4	β_sPU	229
23.3.5	β_vPU	231
23.3.6	β_pPU	231
23.4	Other trees	231
23.4.1	β_{spuARB}	231
23.4.2	$\beta_{spuP_{nx}}$	231
23.4.3	$\beta_{nx}PU_a$	231
V	Appendices	233
A	Future Work	235
A.1	Adjective ordering	235
A.2	More work on Determiners	235
A.3	- <i>ing</i> adjectives	236
A.4	Verb selectional restrictions	236
A.5	Thematic Roles	237
B	Metarules	238
B.1	Introduction	238
B.2	The definition of a metarule in XTAG	239
B.2.1	Node names, variable instantiation, and matches	239
B.2.2	Structural Matching	240
B.2.3	Output Generation	242
B.2.4	Feature Matching	243
B.3	Examples	244

B.4	The Access to the Metarules through the XTAG Interface	246
C	Lexical Organization	250
C.1	Introduction	250
C.2	System Overview	250
C.2.1	Subcategorization frames	251
C.2.2	Blocks	251
C.2.3	Lexical Redistribution Rules (LRRs)	253
C.2.4	Tree generation	253
C.3	Implementation	254
C.4	Generating grammars	254
C.5	Summary	257
D	Tree Naming conventions	258
D.1	Tree Families	258
D.2	Trees within tree families	258
D.3	Assorted Initial Trees	259
D.4	Assorted Auxiliary Trees	259
D.4.1	Relative Clause Trees	260
E	Features	261
E.1	Agreement	261
E.1.1	Agreement and Movement	263
E.2	Case	263
E.2.1	ECM	264
E.2.2	Agreement and Case	264
E.3	Extraction and Inversion	264
E.3.1	Inversion, Part 1	265
E.3.2	Inversion, Part 2	266
E.4	Clause Type	266
E.4.1	Auxiliary Selection	267
E.5	Relative Clauses	267
E.6	Complementizer Selection	269
E.6.1	Verbs with object sentential complements	269
E.6.2	Verbs with sentential subjects	270
E.6.3	<i>That</i> -trace and <i>for</i> -trace effects	271
E.7	Determiner ordering	271
E.8	Punctuation	271
E.9	Conjunction	272
E.10	Comparatives	272
E.11	Control	272
E.12	Other Features	272

F	Evaluation and Results	274
F.1	Parsing Corpora	274
F.2	TSNLP	274
F.3	Chunking and Dependencies in XTAG Derivations	276
F.4	Comparison with IBM	278
F.5	Comparison with Alvey	279
F.6	Comparison with CLARE	279

List of Figures

2.1	Elementary trees in TAG	6
2.2	Substitution in TAG	6
2.3	Adjunction in TAG	7
2.4	Lexicalized Elementary trees	7
2.5	Substitution in FB-LTAG	8
2.6	Adjunction in FB-LTAG	9
2.7	Lexicalized Elementary Trees with Features	10
3.1	XTAG system diagram	12
3.2	Output Structures from the Parser	15
3.3	Interfaces database	16
3.4	XTAG Interface	17
4.1	Different subcategorization frames for the verb <i>buy</i>	20
4.2	Trees illustrating the difference between Complements and Adjuncts	21
4.3	Lexicalized NP trees with case markings	23
4.4	Assigning case in prepositional phrases	24
4.5	Case assignment to NP arguments	24
4.6	Assigning case according to verb form	25
4.7	Proper case assignment with auxiliary verbs	26
6.1	Declarative Intransitive Tree: $\alpha_{nx}0V$	39
6.2	Declarative Transitive Tree: $\alpha_{nx}0V_{nx}1$	39
6.3	Declarative Ditransitive Tree: $\alpha_{nx}0V_{nx}1_{nx}2$	40
6.4	Declarative Ditransitive with PP Tree: $\alpha_{nx}0V_{nx}1_{pnx}2$	41
6.5	Declarative Ditransitive with PP shift Trees: $\alpha_{nx}0V_{nx}1_{Pnx}2$ (a) and $\alpha_{nx}0V_{nx}2_{nx}1$ (b)	42
6.6	Declarative Sentential Complement with NP Tree: $\beta_{nx}0V_{nx}1s2$	44
6.7	Declarative Intransitive Verb Particle Tree: $\alpha_{nx}0V_{pl}$	45
6.8	Declarative Transitive Verb Particle Tree: $\alpha_{nx}0V_{plnx}1$ (a) and $\alpha_{nx}0V_{nx}1_{pl}$ (b)	45
6.9	Declarative Ditransitive Verb Particle Tree: $\alpha_{nx}0V_{plnx}1_{nx}2$ (a), $\alpha_{nx}0V_{nx}1_{plnx}2$ (b) and $\alpha_{nx}0V_{nx}1_{nx}2_{pl}$ (c)	46
6.10	Declarative Intransitive with PP Tree: $\alpha_{nx}0V_{pnx}1$	47
6.11	Declarative PP Complement Tree: $\alpha_{nx}0V_{Pnx}1$	48
6.12	Declarative Sentential Complement Tree: $\beta_{nx}0V_{s1}$	49
6.13	Declarative Intransitive with Adjective Tree: $\alpha_{nx}0V_{ax}1$	50
6.14	Declarative Sentential Subject Tree: $\alpha_{s0}V_{nx}1$	51

6.15	Declarative Light Verb Tree: $\alpha nx0lVN1$	52
6.16	Declarative Light Verbs with PP Tree: $\alpha nx0lVN1Pnx2$ (a), $\alpha nx0lVnx2N1$ (b) . .	52
6.17	Declarative NP It-Cleft Tree: $\alpha ItVpnx1s2$	53
6.18	Declarative PP It-Cleft Tree: $\alpha ItVnx1s2$	54
6.19	Declarative Adverb It-Cleft Tree: $\alpha ItVad1s2$	55
6.20	Declarative Adjective Small Clause Tree: $\alpha nx0Ax1$	55
6.21	Declarative Adjective Small Clause with Sentential Complement Tree: $\alpha nx0A1s1$	56
6.22	Declarative Adjective Small Clause with Sentential Subject Tree: $\alpha s0Ax1$	57
6.23	Declarative Equative <i>BE</i> Tree: $\alpha nx0BEnx1$	57
6.24	Declarative NP Small Clause Trees: $\alpha nx0N1$	58
6.25	Declarative NP with Sentential Complement Small Clause Tree: $\alpha nx0N1s1$. . .	59
6.26	Declarative NP Small Clause with Sentential Subject Tree: $\alpha s0N1$	59
6.27	Declarative PP Small Clause Tree: $\alpha nx0Pnx1$	60
6.28	Declarative Exhaustive PP Small Clause Tree: $\alpha nx0Px1$	61
6.29	Declarative PP Small Clause with Sentential Subject Tree: $\alpha s0Pnx1$	61
6.30	Declarative Intransitive Sentential Subject Tree: $\alpha s0V$	62
6.31	Sentential Subject Tree with ‘to’ complement: $\alpha s0Vtonx1$	62
6.32	Declarative PP Small Clause tree with two-word preposition, where the first word is an adverb, and the second word is a preposition: $\alpha nx0ARBPNx1$	63
6.33	Declarative PP Small Clause tree with two-word preposition, where the first word is an adjective, and the second word is a preposition: $\alpha nx0APnx1$	64
6.34	Declarative PP Small Clause tree with two-word preposition, where the first word is a noun, and the second word is a preposition: $\alpha nx0NPNx1$	64
6.35	Declarative PP Small Clause tree with two-word preposition, where both words are prepositions: $\alpha nx0PPnx1$	65
6.36	Declarative PP Small Clause tree with three-word preposition, where the middle noun is marked for null adjunction: $\alpha nx0PNaPNx1$	66
6.37	Declarative PP Small Clause with Sentential Subject Tree, with two-word prepo- sition, where the first word is an adverb, and the second word is a preposition: $\alpha s0ARBPNx1$	67
6.38	Declarative PP Small Clause with Sentential Subject Tree, with two-word prepo- sition, where the first word is an adjective, and the second word is a preposition: $\alpha s0APnx1$	68
6.39	Declarative PP Small Clause with Sentential Subject Tree, with two-word prepo- sition, where the first word is a noun, and the second word is a preposition: $\alpha s0NPNx1$	68
6.40	Declarative PP Small Clause with Sentential Subject Tree, with two-word prepo- sition, where both words are prepositions: $\alpha s0PPnx1$	69
6.41	Declarative PP Small Clause with Sentential Subject Tree, with three-word preposition, where the middle noun is marked for null adjunction: $\alpha s0PNaPNx1$.	70
6.42	Predicative Adjective with Sentential Subject and Complement: $\alpha s0A1s1$	71
6.43	Declarative Locative Adverbial Small Clause Tree: $\alpha nx0nx1ARB$	71
6.44	Wh-moved Locative Small Clause Tree: $\alpha W1nx0nx1ARB$	72
6.45	ECM Tree: $\beta Xnx0Vs1$	72
6.46	Declarative Transitive Idiom Tree: $\alpha nx0VDN1$	73

6.47	Declarative Idiom with V, D, A, and N Anchors Tree: $\alpha_{nx0VDAN1}$	74
6.48	Declarative Idiom with V and N Anchors Tree: α_{nx0VN1}	75
6.49	Declarative Idiom with V, A, and N Anchors Tree: $\alpha_{nx0VAN1}$	75
6.50	Declarative Idiom with V, D, A, N, and Prep Anchors Tree: $\alpha_{nx0VDAN1P_{nx2}}$.	76
6.51	Declarative Idiom with V, A, N, and Prep Anchors Tree: $\alpha_{nx0VAN1P_{nx2}}$	77
6.52	Declarative Idiom with V, N, and Prep Anchors Tree: $\alpha_{nx0VN1P_{nx2}}$	77
6.53	Declarative Idiom with V, D, N, and Prep Anchors Tree: $\alpha_{nx0VDN1P_{nx2}}$	78
7.1	Ergative Tree: α_{Enx1V}	80
8.1	Tree β COMPs, anchored by <i>that</i>	85
8.2	Sentential complement tree: β_{nx0Vs1}	86
8.3	Trees for <i>The emu thinks that the aardvark smells terrible.</i>	87
8.4	Tree for <i>Who smells terrible?</i>	87
8.5	ECM tree: $\beta_{Xnx0Vs1}$	89
8.6	Sample ECM parse	89
8.7	ECM passive	90
8.8	Comparison of < assign-comp > values for sentential subjects: $\alpha_{s0V_{nx1}}$ (a) and sentential complements: β_{nx0Vs1} (b)	91
8.9	Sample trees for preposition: β_{Pss} (a) and noun: α_{NXNs} (b) taking sentential complements	92
8.10	Tree for <i>persuaded</i>	94
8.11	Tree for <i>leave</i>	94
8.12	Derived tree for <i>Srini persuaded Mickey to leave</i>	95
8.13	Derivation tree for <i>Srini persuaded Mickey to leave</i>	95
9.1	Predicative trees: α_{nx0N1} (a), α_{nx0Ax1} (b) and $\alpha_{nx0P_{nx1}}$ (c)	101
9.2	Copula auxiliary tree: $\beta_{V_{vx}}$	102
9.3	Predicative AP tree with features: α_{nx0Ax1}	104
9.4	<i>Consider</i> tree for embedded small clauses	105
9.5	Raising verb with experiencer tree: $\beta_{V_{pxvx}}$	105
9.6	Raising adjective tree: $\beta_{V_{vx}-adj}$	105
9.7	Equative <i>BE</i> trees: $\alpha_{nx0B_{Enx1}}$ (a) and $\alpha_{Inv_{nx0}B_{Enx1}}$ (b)	106
10.1	Dative shift trees: $\alpha_{nx0V_{nx1}P_{nx2}}$ (a) and $\alpha_{nx0V_{nx2}nx1}$ (b)	108
11.1	It-cleft with PP clefted element: $\alpha_{ItV_{pnx1s2}}$ (a) and $\alpha_{InvItV_{pnx1s2}}$ (b)	111
12.1	Passive trees in the Sentential Complement with NP tree family: β_{nx1Vs2} (a), $\beta_{nx1V_{bynx0s2}}$ (b) and $\beta_{nx1Vs2bynx0}$ (c)	115
13.1	Transitive tree with object extraction: $\alpha_{W1nx0V_{nx1}}$	118
13.2	Intransitive tree with subject extraction: α_{W0nx0V}	121
13.3	Ditransitive trees with direct object: $\alpha_{W1nx0V_{nx1nx2}}$ (a) and indirect object extraction: $\alpha_{W2nx0V_{nx1nx2}}$ (b)	122
13.4	Ditransitive with PP tree with the object of the PP extracted: $\alpha_{W2nx0V_{nx1pnx2}}$ 122	
13.5	Ditransitive with PP with PP extraction tree: $\alpha_{pW2nx0V_{nx1pnx2}}$	123

13.6	Predicative Adjective tree with extracted adjective: $\alpha\text{WA1nx0Vax1}$	124
14.1	Relative clause trees in the transitive tree family: $\beta\text{Nc1nx0Vnx1}$ (a) and $\beta\text{N0nx0Vnx1}$ (b)	126
14.2	Adjunct relative clause tree with PP-pied-piping in the transitive tree family: $\beta\text{Npxnx0Vnx1}$	128
14.3	Determiner tree with <rel-clause> feature: βDnx	132
15.1	Auxiliary Trees for Subordinating Conjunctions	136
15.2	Trees Anchored by Subordinating Conjunctions: $\beta\text{vxPARBPs}$ and $\beta\text{vxParbPs}$	138
15.3	Sample Participial Adjuncts	139
15.4	Example of discourse conjunction, from Seuss' <i>The Lorax</i>	139
16.1	Transitive imperative tree: $\alpha\text{Inx0Vnx1}$	141
17.1	Determiner Gerund tree from the transitive tree family: $\alpha\text{Dnx0Vnx1}$	144
17.2	NP Gerund tree from the transitive tree family: $\alpha\text{Gnx0Vnx1}$	145
17.3	Passive Gerund trees from the transitive tree family: $\alpha\text{Gnx1Vbynx0}$ (a) and αGnx1V (b)	147
18.1	NP Tree	152
18.2	Determiner Trees with Features	153
18.3	Multi-word Determiner tree: βDDnx	158
18.4	Genitive Determiner Tree	160
18.5	Genitive NP tree for substitution: αDnxG	160
18.6	Partitive Determiner Tree	162
18.7	(a) Adverb modifying a determiner; (b) Adverb modifying a noun phrase	163
19.1	Standard Tree for Adjective modifying a Noun: βAn	166
19.2	Multiple adjectives modifying a noun	167
19.3	Noun-noun compounding tree: βNn (not all features displayed)	168
19.4	Time Phrase Modifier trees: βNs , βNvx , βvxN , βnxN	170
19.5	Time NPs with and without a determiner	170
19.6	Time NP trees: Two different attachments	171
19.7	Time NPs in different positions (βvxN , βnxN and βNs)	171
19.8	Time NPs: Derived tree and Derivation (βNvx position)	172
19.9	Selected Prepositional Phrase Modifier trees: βPss , βnxPnx , βvxP and βvxPPnx	173
19.10	Adverb Trees for pre-modification of S: βARBs (a) and post-modification of a VP: βvxARB (b)	174
19.11	Derived tree for <i>How did you fall?</i>	176
19.12	Complex adverb phrase modifier: $\beta\text{ARBarbs}$	177
19.13	Selected Focus and Multi-word Adverb Modifier trees: βARBnx , βPARBd and βPaPd	178
19.14	Selected Multi-word Adverb Modifier trees (for adverbs like <i>sort of</i> , <i>kind of</i>): βNPax , βNPvx , βvxNP	179
19.15	Selected Multi-word Adverb Modifier trees (for adverbs like <i>a little</i> , <i>a bit</i>): βvxDA , βDAax , βDNpx	180

19.16	Locative Modifier Trees: β_{nxnxARB} , β_{nxARB}	180
19.17	Locative Phrases featuring NP and Adverb Degree Specifications	181
20.1	Auxiliary verb tree for non-inverted sentences: βV_{vx}	185
20.2	Auxiliary trees for <i>The music should have been being played</i>	186
20.3	<i>The music should have been being played</i>	187
20.4	Trees for auxiliary verb inversion: βV_{s} (a) and βV_{vx} (b)	188
20.5	<i>will John buy a backpack ?</i>	189
21.1	Tree for adjective conjunction: $\beta a1\text{CONJa}2$ and a resulting parse tree	195
21.2	Tree for NP conjunction: $\beta \text{CONJnx}1\text{CONJnx}2$ and a resulting parse tree	196
21.3	Tree for determiner conjunction: $\beta d1\text{CONJd}2.\text{ps}$	197
21.4	Tree for sentential conjunction: $\beta s1\text{CONJs}2$	198
21.5	198
21.6	$\beta a1\text{CONJa}2$ (a) anchored by comma and (b) anchored by <i>and</i>	199
21.7	Tree for conjunction with but-not: $\beta \text{px}1\text{CONJARBpx}2$	200
21.8	Tree for conjunction with not-but: $\beta \text{ARBnx}1\text{CONJnx}2$	201
21.9	Example of conjunction with <i>to</i>	202
21.10	Coordination schema	202
21.11	An example of the <i>conjoin</i> operation. $\{1\}$ denotes a shared dependency.	203
21.12	Coordination as adjunction.	203
22.1	Tree for Metalinguistic Adjective Comparative: βARBaPa	206
22.2	Tree for Adjective-Extreme Comparative: βARBPa	207
22.3	Nominal comparative trees	209
22.4	Tree for Lone Comparatives: αCARB	210
22.5	Comparative conjunctions.	211
22.6	Comparative conjunctions.	212
22.7	Adjunction of β_{nxPnx} to NP modified by comparative adjective.	213
22.8	Elliptical adjectival comparative trees	214
22.9	Comparativized adjective triggering βCnxPnx	215
22.10	Adjunction of β_{axPnx} to comparative adjective.	216
22.11	Adverbial comparative trees	217
23.1	The β_{nxPUnxPU} tree, anchored by parentheses	221
23.2	An N-level modifier, using the β_{nPUnx} tree	222
23.3	The derived trees for an NP with (a) a peripheral, dash-separated appositive and (b) an NP colon expansion (uttered by the Mouse in <i>Alice's Adventures in Wonderland</i>)	223
23.4	The $\beta \text{PUpxPUvx}$ tree, anchored by commas	223
23.5	Tree illustrating the use of $\beta \text{PUpxPUvx}$	224
23.6	A tree illustrating the use of sPUnx for a colon expansion attached at S.	225
23.7	$\beta \text{PU}s\text{PU}$ anchored by parentheses, and in a derivation, along with βPUnxPU	226
23.8	$\beta \text{PU}s$, with features displayed	228
23.9	$\beta \text{sPU}s$, with features displayed	230

B.1	Metarule for wh-movement of subject	245
B.2	Metarule for wh-movement of object	245
B.3	Metarule for general wh movement of an NP	246
B.4	Application of wh-movement rule to Tnx0Vnx1Pnx2	247
B.5	Parallel application of metarules	248
B.6	Sequential application of metarules	248
B.7	Cumulative application of metarules	249
C.1	Lexical Organization: System Overview	251
C.2	Some subcategorization blocks	252
C.3	Transformation blocks for extraction	252
C.4	Elementary trees generated from combining blocks	253
C.5	Partial inheritance lattice in English	254
C.6	Implementation of the system	255
C.7	Interface for creating a grammar	255
C.8	Part of the Interface for creating blocks	256

Abstract

This document describes a sizable grammar of English written in the TAG formalism and implemented for use with the XTAG system. This report and the grammar described herein supersedes the TAG grammar described in [XTAG-Group, 1995]. The English grammar described in this report is based on the TAG formalism developed in [Joshi *et al.*, 1975], which has been extended to include lexicalization ([Schabes *et al.*, 1988]), and unification-based feature structures ([Vijay-Shanker and Joshi, 1991]). The range of syntactic phenomena that can be handled is large and includes auxiliaries (including inversion), copula, raising and small clause constructions, topicalization, relative clauses, infinitives, gerunds, passives, adjuncts, it-clefts, wh-clefts, PRO constructions, noun-noun modifications, extraposition, determiner sequences, genitives, negation, noun-verb contractions, sentential adjuncts and imperatives. This technical report corresponds to the XTAG Release 8/31/98. The XTAG grammar is continuously updated with the addition of new analyses and modification of old ones, and an online version of this report can be found at the XTAG web page: <http://www.cis.upenn.edu/~xtag/>.

Acknowledgements

We are immensely grateful to Aravind Joshi for supporting this project.

The following people have contributed to the development of grammars in the project: Anne Abeille, Jason Baldridge, Rajesh Bhatt, Kathleen Bishop, Raman Chandrasekar, Sharon Cote, Beatrice Daille, Christine Doran, Dania Egedi, Tim Farrington, Jason Frank, Caroline Heycock, Beth Ann Hockey, Roumyana Izvorski, Karin Kipper, Daniel Karp, Seth Kulick, Young-Suk Lee, Heather Matayek, Patrick Martin, Megan Moser, Sabine Petillon, Rashmi Prasad, Laura Siegel, Yves Schabes, Victoria Tredinnick and Raffaella Zanuttini.

The XTAG system has been developed by: Tilman Becker, Richard Billington, Andrew Chalnack, Dania Egedi, Devtosh Khare, Albert Lee, David Magerman, Alex Mallet, Patrick Paroubek, Rich Pito, Gilles Prigent, Carlos Prolo, Anoop Sarkar, Yves Schabes, William Schuler, B. Srinivas, Fei Xia, Yuji Yoshiie and Martin Zaidel.

We would also like to thank Michael Hegarty, Lauri Karttunen, Anthony Kroch, Mitchell Marcus, Martha Palmer, Owen Rambow, Philip Resnik, Beatrice Santorini and Mark Steedman.

In addition, Jeff Aaronson, Douglas DeCarlo, Mark-Jason Dominus, Mark Foster, Gaylord Holder, David Magerman, Ken Noble, Steven Shapiro and Ira Winston have provided technical support. Administrative support was provided by Susan Deysher, Carolyn Elken, Jodi Kerper, Christine Sandy and Trisha Yannuzzi.

This work was partially supported by NSF Grant SBR8920230 and ARO Grant DAAH0404-94-G-0426.

Part I

General Information

Chapter 1

Getting Around

This technical report presents the English XTAG grammar as implemented by the XTAG Research Group at the University of Pennsylvania. The technical report is organized into four parts, plus a set of appendices. Part 1 contains general information about the XTAG system and some of the underlying mechanisms that help shape the grammar. Chapter 2 contains an introduction to the formalism behind the grammar and parser, while Chapter 3 contains information about the entire XTAG system. Linguists interested solely in the grammar of the XTAG system may safely skip Chapters 2 and 3. Chapter 4 contains information on some of the linguistic principles that underlie the XTAG grammar, including the distinction between complements and adjuncts, and how case is handled.

The actual description of the grammar begins with Part 2, and is contained in the following three parts. Parts 2 and 3 contains information on the verb classes and the types of trees allowed within the verb classes, respectively, while Part 4 contains information on trees not included in the verb classes (e.g. NP's, PP's, various modifiers, etc). Chapter 5 of Part 2 contains a table that attempts to provide an overview of the verb classes and tree types by providing a graphical indication of which tree types are allowed in which verb classes. This has been cross-indexed to tree figures shown in the tech report. Chapter 6 contains an overview of all of the verb classes in the XTAG grammar. The rest of Part 2 contains more details on several of the more interesting verb classes, including ergatives, sentential subjects, sentential complements, small classes, ditransitives, and it-clefts.

Part 3 contains information on some of the tree types that are available within the verb classes. These tree types correspond to what would be transformations in a movement based approach. Not all of these types of trees are contained in all of the verb classes. The table (previously mentioned) in Part 2 contains a list of the tree types and indicates which verb classes each occurs in.

Part 4 focuses on the non-verb class trees in the grammar. NP's and determiners are presented in Chapter 18, while the various modifier trees are presented in Chapter 19. Auxiliary verbs, which are classed separate from the verb classes, are presented in Chapter 20, while certain types of conjunction are shown in Chapter 21. The XTAG treatment of comparatives is presented in Chapter 22, and our treatment of punctuation is discussed in Chapter 23.

Throughout the technical report, mention is occasionally made of changes or analyses that we hope to incorporate in the future. Appendix A details a list of these and other future work. The appendices also contain information on some of the nitty gritty details of the

CHAPTER 1. GETTING AROUND

XTAG grammar, including a system of metarules which can be used for grammar development and maintenance in Appendix B, a system for the organization of the grammar in terms of an inheritance hierarchy is in Appendix C, the tree naming conventions used in XTAG are explained in detail in Appendix D, and a comprehensive list of the features used in the grammar is given in Appendix E. Appendix F contains an evaluation of the XTAG grammar, including comparisons with other wide coverage grammars.

Chapter 2

Feature-Based, Lexicalized Tree Adjoining Grammars

The English grammar described in this report is based on the TAG formalism ([Joshi *et al.*, 1975]), which has been extended to include lexicalization ([Schabes *et al.*, 1988]), and unification-based feature structures ([Vijay-Shanker and Joshi, 1991]). Tree Adjoining Languages (TALs) fall into the class of mildly context-sensitive languages, and as such are more powerful than context free languages. The TAG formalism in general, and lexicalized TAGs in particular, are well-suited for linguistic applications. As first shown by [Joshi, 1985] and [Kroch and Joshi, 1987], the properties of TAGs permit us to encapsulate diverse syntactic phenomena in a very natural way. For example, TAG's extended domain of locality and its factoring of recursion from local dependencies lead, among other things, to a localization of so-called unbounded dependencies.

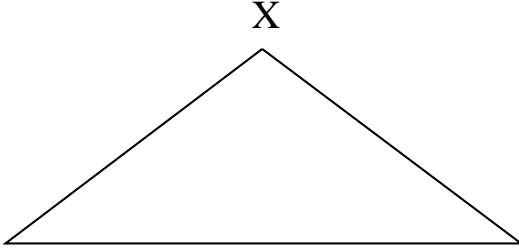
2.1 TAG formalism

The primitive elements of the standard TAG formalism are known as elementary trees. ELEMENTARY TREES are of two types: initial trees and auxiliary trees (see Figure 2.1). In describing natural language, INITIAL TREES are minimal linguistic structures that contain no recursion, i.e. trees containing the phrasal structure of simple sentences, NP's, PP's, and so forth. Initial trees are characterized by the following: 1) all internal nodes are labeled by non-terminals, 2) all leaf nodes are labeled by terminals, or by non-terminal nodes marked for substitution. An initial tree is called an X-type initial tree if its root is labeled with type X.

Recursive structures are represented by AUXILIARY TREES, which represent constituents that are adjuncts to basic structures (e.g. adverbials). Auxiliary trees are characterized as follows: 1) all internal nodes are labeled by non-terminals, 2) all leaf nodes are labeled by terminals, or by non-terminal nodes marked for substitution, except for exactly one non-terminal node, called the foot node, which can only be used to adjoin the tree to another node¹, 3) the foot node has the same label as the root node of the tree.

¹A null adjunction constraint (NA) is systematically put on the foot node of an auxiliary tree. This disallows adjunction of a tree onto the foot node itself.

Initial Tree:



Auxiliary Tree:

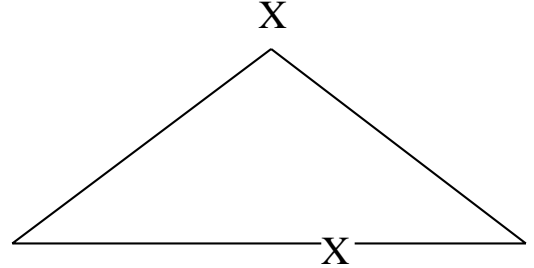


Figure 2.1: Elementary trees in TAG

There are two operations defined in the TAG formalism, substitution² and adjunction. In the SUBSTITUTION operation, the root node on an initial tree is merged into a non-terminal leaf node marked for substitution in another initial tree, producing a new tree. The root node and the substitution node must have the same name. Figure 2.2 shows two initial trees and the tree resulting from the substitution of one tree into the other.

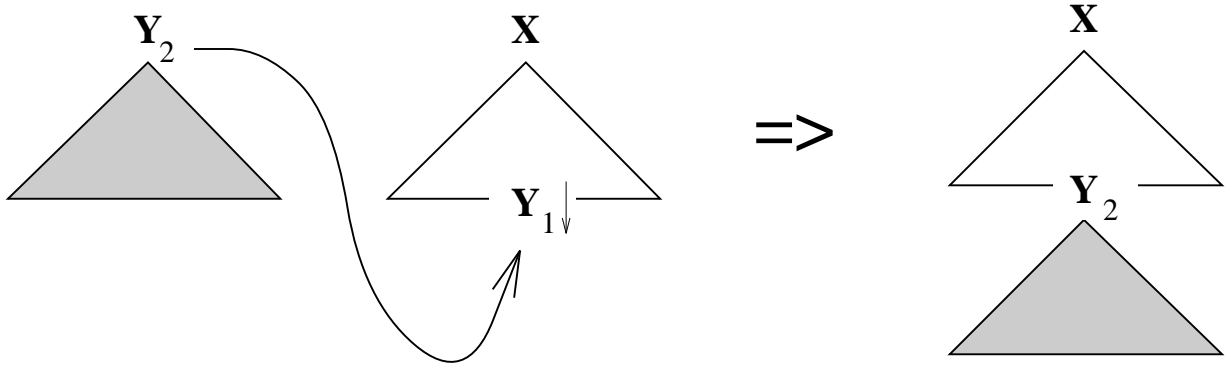


Figure 2.2: Substitution in TAG

In an ADJUNCTION operation, an auxiliary tree is grafted onto a non-terminal node anywhere in an initial tree. The root and foot nodes of the auxiliary tree must match the node at which the auxiliary tree adjoins. Figure 2.3 shows an auxiliary tree and an initial tree, and the tree resulting from an adjunction operation.

A TAG G is a collection of finite initial trees, I , and auxiliary trees, A . The TREE SET of a TAG G , $\mathcal{T}(G)$ is defined to be the set of all derived trees starting from S-type initial trees in I whose frontier consists of terminal nodes (all substitution nodes having been filled). The STRING LANGUAGE generated by a TAG, $\mathcal{L}(G)$, is defined to be the set of all terminal strings on the frontier of the trees in $\mathcal{T}(G)$.

²Technically, substitution is a specialized version of adjunction, but it is useful to make a distinction between the two.

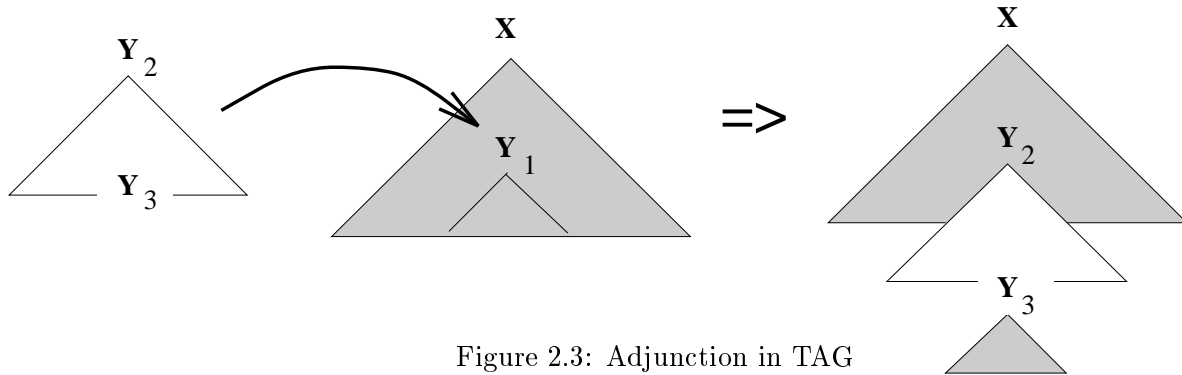


Figure 2.3: Adjunction in TAG

2.2 Lexicalization

‘Lexicalized’ grammars systematically associate each elementary structure with a lexical anchor. This means that in each structure there is a lexical item that is realized. It does not mean simply adding feature structures (such as head) and unification equations to the rules of the formalism. These resultant elementary structures specify extended domains of locality (as compared to CFGs) over which constraints can be stated.

Following [Schabes *et al.*, 1988] we say that a grammar is **LEXICALIZED** if it consists of 1) a finite set of structures each associated with a lexical item, and 2) an operation or operations for composing the structures. Each lexical item will be called the **ANCHOR** of the corresponding structure, which defines the domain of locality over which constraints are specified. Note then, that constraints are local with respect to their anchor.

Not every grammar is in a lexicalized form.³ In the process of lexicalizing a grammar, the lexicalized grammar is required to be strongly equivalent to the original grammar, i.e. it must produce not only the same language, but the same structures or tree set as well.

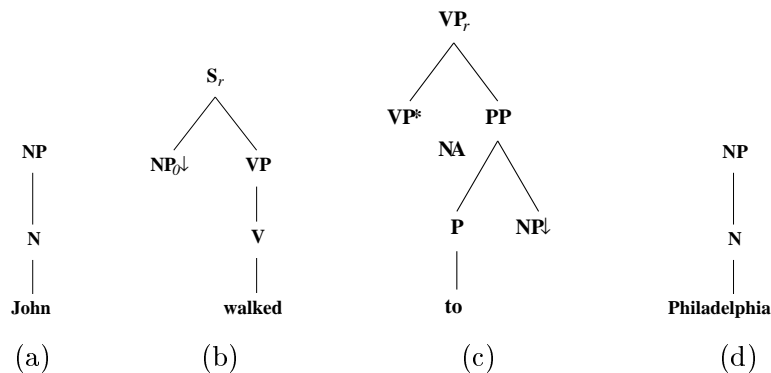


Figure 2.4: Lexicalized Elementary trees

³Notice the similarity of the definition of a lexicalized grammar with the off line parsability constraint ([Kaplan and Bresnan, 1983]). As consequences of our definition, each structure has at least one lexical item (its anchor) attached to it and all sentences are finitely ambiguous.

In Figure 2.4, which shows sample initial and auxiliary trees, substitution sites are marked by a \downarrow , and foot nodes are marked by an $*$. This notation is standard and is followed in the rest of this report.

2.3 Unification-based features

In a unification framework, a feature structure is associated with each node in an elementary tree. This feature structure contains information about how the node interacts with other nodes in the tree. It consists of a top part, which generally contains information relating to the supernode, and a bottom part, which generally contains information relating to the subnode. Substitution nodes, however, have only the top features, since the tree substituting in logically carries the bottom features.

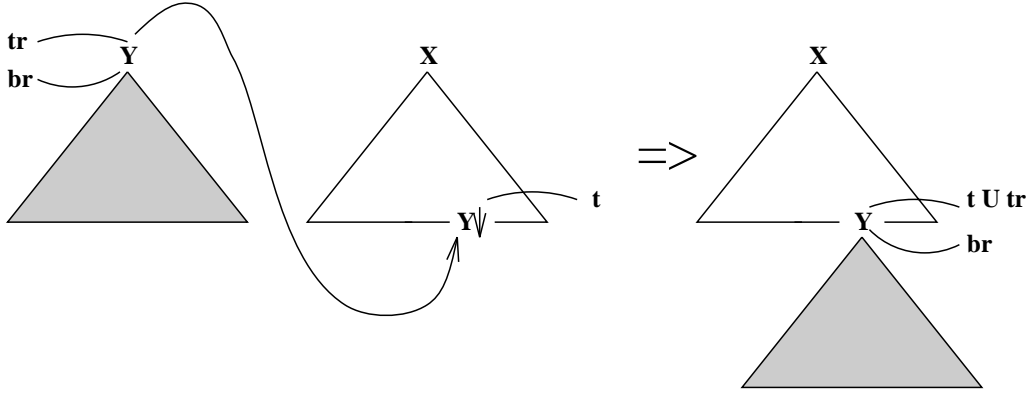


Figure 2.5: Substitution in FB-LTAG

The notions of substitution and adjunction must be augmented to fit within this new framework. The feature structure of a new node created by substitution inherits the union of the features of the original nodes. The top feature of the new node is the union of the top features of the two original nodes, while the bottom feature of the new node is simply the bottom feature of the top node of the substituting tree (since the substitution node has no bottom feature). Figure 2.5⁴ shows this more clearly.

Adjunction is only slightly more complicated. The node being adjoined into splits, and its top feature unifies with the top feature of the root adjoining node, while its bottom feature unifies with the bottom feature of the foot adjoining node. Again, this is easier shown graphically, as in Figure 2.6⁵.

The embedding of the TAG formalism in a unification framework allows us to dynamically specify local constraints that would have otherwise had to have been made statically within the trees. Constraints that verbs make on their complements, for instance, can be implemented through the feature structures. The notions of Obligatory and Selective Adjunction, crucial

⁴abbreviations in the figure: t=top feature structure, tr=top feature structure of the root, br=bottom feature structure of the root, U=unification

⁵abbreviations in the figure: t=top feature structure, b=bottom feature structure, tr=top feature structure of the root, br=bottom feature structure of the root, tf=top feature structure of the foot, bf=bottom feature structure of the foot, U=unification

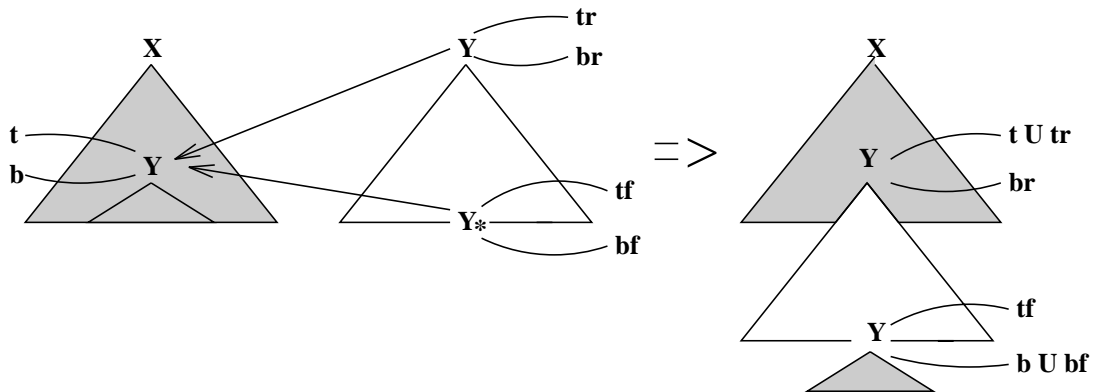


Figure 2.6: Adjunction in FB-LTAG

to the formation of lexicalized grammars, can also be handled through the use of features.⁶ Perhaps more important to developing a grammar, though, is that the trees can serve as a schemata to be instantiated with lexical-specific features when an anchor is associated with the tree. To illustrate this, Figure 2.7 shows the same tree lexicalized with two different verbs, each of which instantiates the features of the tree according to its lexical selectional restrictions.

In Figure 2.7, the lexical item *thinks* takes an indicative sentential complement, as in the sentence *John thinks that Mary loves Sally*. *Want* takes a sentential complement as well, but an infinitive one, as in *John wants to love Mary*. This distinction is easily captured in the features and passed to other nodes to constrain which trees this tree can adjoin into, both cutting down the number of separate trees needed and enforcing conceptual Selective Adjunctions (SA).

⁶The remaining constraint, Null Adjunction (NA), must still be specified directly on a node.

CHAPTER 2. FEATURE-BASED, LEXICALIZED TREE ADJOINING GRAMMARS

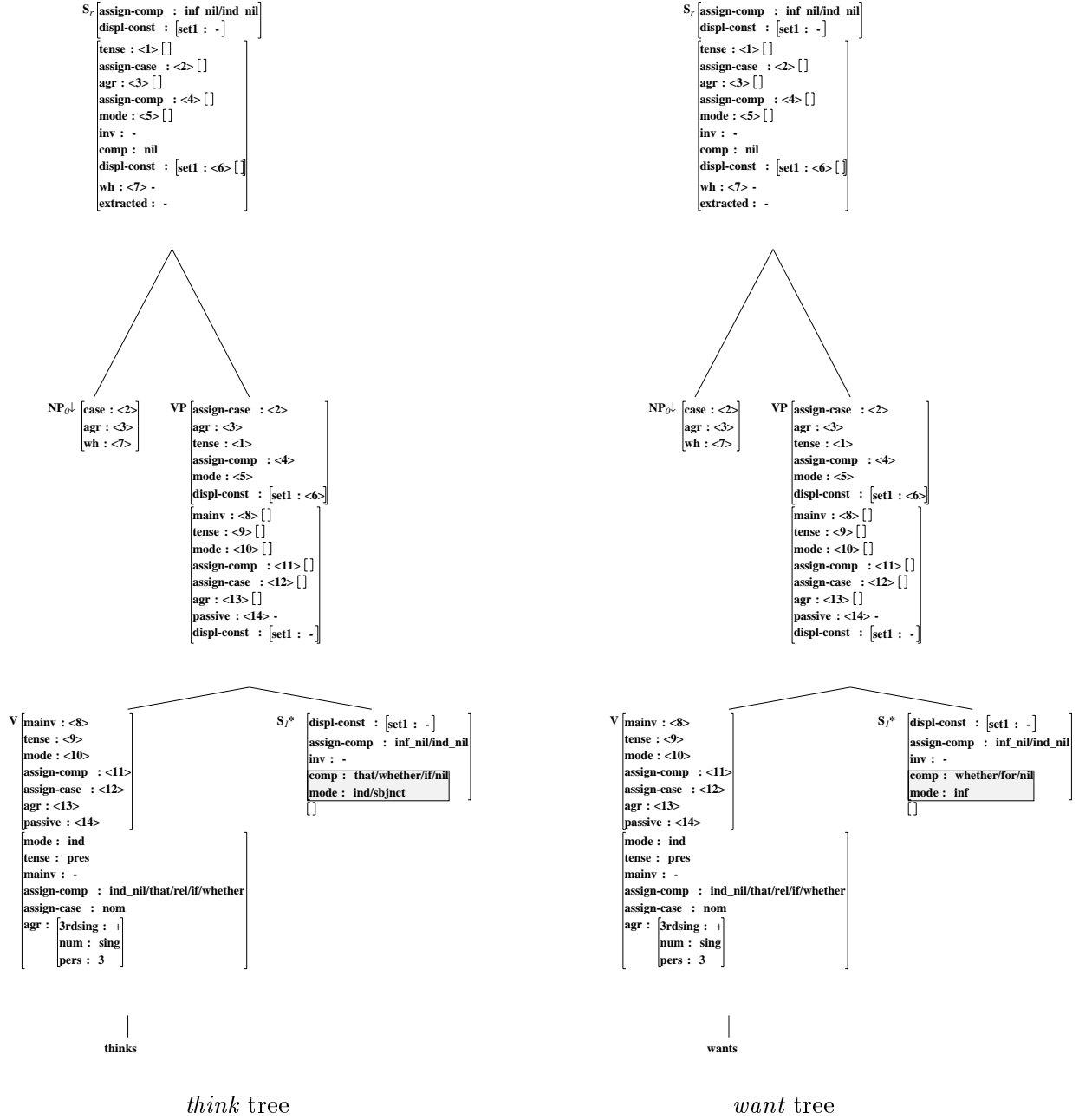


Figure 2.7: Lexicalized Elementary Trees with Features

Chapter 3

Overview of the XTAG System

This section focuses on the various components that comprise the parser and English grammar in the XTAG system. Persons interested only in the linguistic analyses in the grammar may skip this section without loss of continuity, although a quick glance at the tagset used in XTAG and the set of non-terminal labels used will be useful. We may occasionally refer back to the various components mentioned in this section.

3.1 System Description

Figure 3.1 shows the overall flow of the system when parsing a sentence; a summary of each component is presented in Table 3.1. At the heart of the system is a parser for lexicalized TAGs ([Schabes and Joshi, 1988; Schabes, 1990]) which produces all legitimate parses for the sentence. The parser has two phases: **Tree Selection** and **Tree Grafting**.

3.1.1 Tree Selection

Since we are working with lexicalized TAGs, each word in the sentence selects at least one tree. The advantage of a lexicalized formalism like LTAGs is that rather than parsing with all the trees in the grammar, we can parse with only the trees selected by the words in the input sentence.

In the XTAG system, the selection of trees by the words is done in several steps. Each step attempts to reduce ambiguity, i.e. reduce the number of trees selected by the words in the sentence.

Morphological Analysis and POS Tagging The input sentence is first submitted to the **Morphological Analyzer** and the **Tagger**. The morphological analyzer ([Karp *et al.*, 1992]) consists of a disk-based database (a compiled version of the derivational rules) which is used to map an inflected word into its stem, part of speech and feature equations corresponding to inflectional information. These features are inserted at the anchor node of the tree eventually selected by the stem. The POS Tagger can be disabled in which case only information from the morphological analyzer is used. The morphology data was originally extracted from the Collins English Dictionary ([Hanks, 1979]) and Oxford Advanced Learner's Dictionary ([Hornby, 1974]) available through ACL-DCI ([Lieberman, 1989]), and then cleaned up and augmented by hand ([Karp *et al.*, 1992]).

CHAPTER 3. OVERVIEW OF THE XTAG SYSTEM

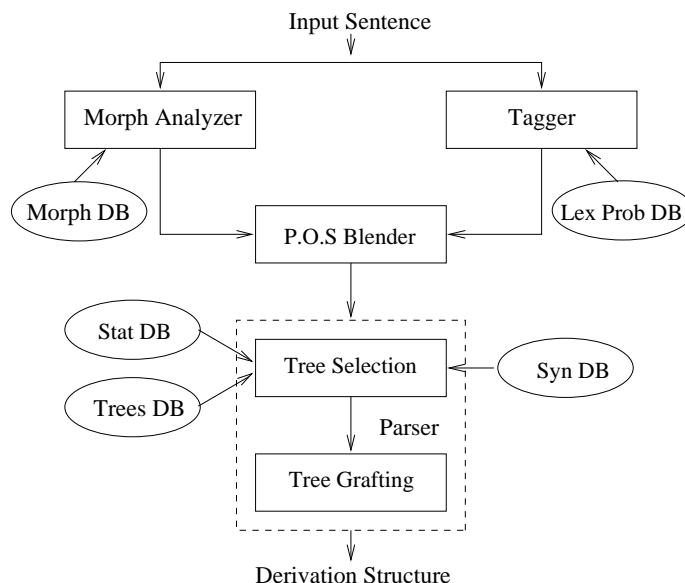


Figure 3.1: Overview of XTAG system

POS Blender The output from the morphological analyzer and the POS tagger go into the **POS Blender** which uses the output of the POS tagger as a filter on the output of the morphological analyzer. Any words that are not found in the morphological database are assigned the POS given by the tagger.

Syntactic Database The syntactic database contains the mapping between particular stem(s) and the tree templates or tree-families stored in the **Tree Database** (see Table 3.1). The syntactic database also contains a list of feature equations that capture lexical idiosyncrasies. The output of the POS Blender is used to search the **Syntactic Database** to produce a set of lexicalized trees with the feature equations associated with the word(s) in the syntactic database unified with the feature equations associated with the trees. Note that the features in the syntactic database can be assigned to any node in the tree and not just to the anchor node. The syntactic database entries were originally extracted from the Oxford Advanced Learner’s Dictionary ([Hornby, 1974]) and Oxford Dictionary for Contemporary Idiomatic English ([Cowie and Mackin, 1975]) available through ACL-DCI ([Lieberman, 1989]), and then modified and augmented by hand ([Egedi and Martin, 1994]). There are more than 31,000 syntactic database entries.¹ Selected entries from this database are shown in Table 3.2.

Default Assignment For words that are not found in the syntactic database, default trees and tree-families are assigned based on their POS tag.

Filters Some of the lexicalized trees chosen in previous stages can be eliminated in order to reduce ambiguity. Two methods are currently used: structural filters which eliminate trees which have impossible spans over the input sentence and a statistical filter based on

¹This number does not include trees assigned by default based on the part-of-speech of the word.

Component	Details
Morphological Analyzer and Morph Database	Consists of approximately 317,000 inflected items derived from over 90000 stems. Entries are indexed on the inflected form and return the root form, POS, and inflectional information.
POS Tagger and Lex Prob Database	Wall Street Journal-trained trigram tagger ([Church, 1988]) extended to output N-best POS sequences ([Soong and Huang, 1990]). Decreases the time to parse a sentence by an average of 93%.
Syntactic Database	More than 30,000 entries. Each entry consists of: the uninflected form of the word, its POS, the list of trees or tree-families associated with the word, and a list of feature equations that capture lexical idiosyncrasies.
Tree Database	1094 trees, divided into 52 tree families and 218 individual trees. Tree families represent subcategorization frames; the trees in a tree family would be related to each other transformationally in a movement-based approach.
X-Interface	Menu-based facility for creating and modifying tree files. User controlled parser parameters: parser's start category, enable/disable/retry on failure for POS tagger. Storage/retrieval facilities for elementary and parsed trees. Graphical displays of tree and feature data structures. Hand combination of trees by adjunction or substitution for grammar development. Ability to manually assign POS tag and/or Supertag before parsing

Table 3.1: System Summary

unigram probabilities of non-lexicalized trees (from a hand corrected set of approximately 6000 parsed sentences). These methods speed the runtime by approximately 87%.

Supertagging Before parsing, one can avail of an optional step of *supertagging* the sentence. This step uses statistical disambiguation to assign a unique elementary tree (or *supertag*) to each word in the sentence. These assignments can then be hand-corrected. These supertags are used as a filter on the tree assignments made so far. More information on supertagging can be found in ([Srinivas, 1997a; Srinivas, 1997b]).

3.1.2 Tree Database

The **Tree Database** contains the tree templates that are lexicalized by following the various steps given above. The lexical items are inserted into distinguished nodes in the tree template called the *anchor nodes*. The part of speech of each word in the sentence corresponds to the label of the anchor node of the trees. Hence the tagset used by the POS Tagger corresponds exactly to the labels of the anchor nodes in the trees. The tagset used in the XTAG system is given in Table 3.3. The tree templates are subdivided into tree families (for verbs and other

```

<<INDEX>>porousness<<ENTRY>>porousness<<POS>>N
<<TREES>>^BNXN ^BN ^CNn
<<FEATURES>>#N_card- #N_const- #N_decreas- #N_definite- #N_gen-
#N_quan- #N_refl-

<<INDEX>>coo<<ENTRY>>coo<<POS>>V<<FAMILY>>TnxOV

<<INDEX>>engross<<ENTRY>>engross<<POS>>V<<FAMILY>>TnxOVnx1
<<FEATURES>>#TRANS+

<<INDEX>>forbear<<ENTRY>>forbear<<POS>>V<<FAMILY>>TnxOVs1
<<FEATURES>>#S1_WH- #S1_inf_for_nil

<<INDEX>>have<<ENTRY>>have<<POS>>V<<ENTRY>>out<<POS>>PL
<<FAMILY>>TnxOVplnx1

```

Table 3.2: Example Syntactic Database Entries.

predicates), and tree files which are simply collections of trees for lexical items like prepositions, determiners, etc².

3.1.3 Tree Grafting

Once a particular set of lexicalized trees for the sentence have been selected, XTAG uses an Earley-style predictive left-to-right parsing algorithm for LTAGs ([Schabes and Joshi, 1988; Schabes, 1990]) to find all derivations for the sentence. The derivation trees and the associated derived trees can be viewed using the X-interface (see Table 3.1). The X-interface can also be used to save particular derivations to disk.

The output of the parser for the sentence *I had a map yesterday* is illustrated in Figure 3.2. The parse tree³ represents the surface constituent structure, while the derivation tree represents the derivation history of the parse. The nodes of the derivation tree are the tree names anchored by the lexical items⁴. The composition operation is indicated by the nature of the arcs: a dashed line is used for substitution and a bold line for adjunction. The number beside each tree name is the address of the node at which the operation took place. The derivation tree can also be interpreted as a dependency graph with unlabeled arcs between words of the sentence.

3.1.4 The Grammar Development Environment

Working with and developing a large grammar is a challenging process, and the importance of having good visualization tools cannot be over-emphasized. Currently the XTAG system has

²The nonterminals in the tree database are A, AP, Ad, AdvP, Comp, Conj, D, N, NP, P, PP, Punct, S, V, VP.

³The feature structures associated with each node of the parse tree are not shown here.

⁴Appendix D explains the conventions used in naming the trees.

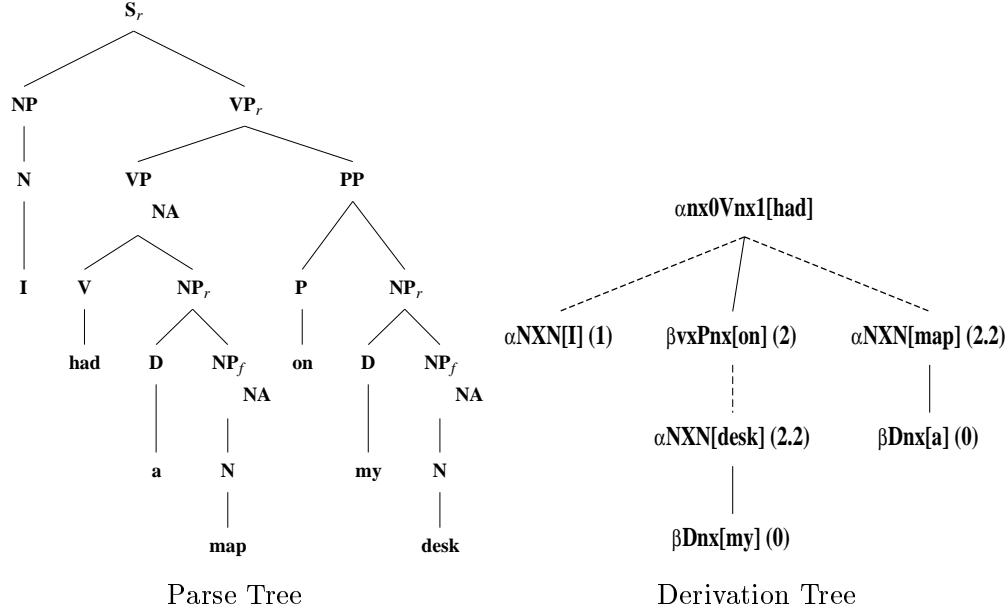


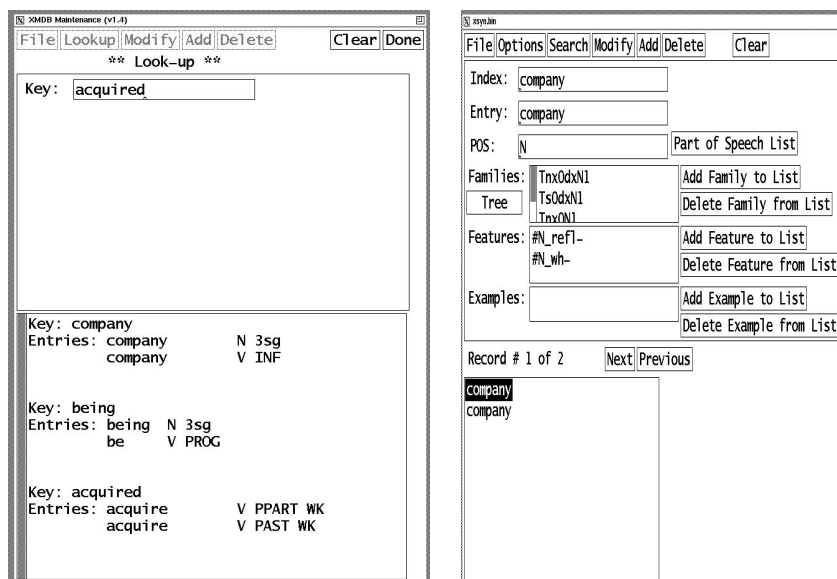
Figure 3.2: Output Structures from the Parser

Part of Speech	Description
A	Adjective
Ad	Adverb
Comp	Complementizer
D	Determiner
G	Genitive Noun
I	Interjection
N	Noun
P	Preposition
PL	Particle
Punct	Punctuation
V	Verb

Table 3.3: XTAG tagset

X-windows based tools for viewing and updating the morphological and syntactic databases ([Karp *et al.*, 1992; Egedi and Martin, 1994]). These are available in both ASCII and binary-encoded database format. The ASCII format is well-suited for various UNIX utilities (awk, sed, grep) while the database format is used for fast access during program execution. However even the ASCII formatted representation is not well-suited for human readability. An X-windows interface for the databases allows users to easily examine them. Searching for specific information on certain fields of the syntactic database is also available. Also, the interface allows a user to insert, delete and update any information in the databases. Figure 3.3(a) shows the interface for the morphology database and Figure 3.3(b) shows the interface for the syntactic database.

XTAG also has a parsing and grammar development interface ([Paroubek *et al.*, 1992]). This



(a) Morphology database

(b) Syntactic database

Figure 3.3: Interfaces to the database maintenance tools

interface includes a tree editor, the ability to vary parameters in the parser, work with multiple grammars and/or parsers, and use metarules for more efficient tree editing and construction ([Becker, 1994]). The interface is shown in Figure 3.4. It has the following features:

- Menu-based facility for creating and modifying tree files and loading grammar files.
- User controlled parser parameters, including the root category (main S, embedded S, NP, etc.), and the use of the tagger (on/off/retry on failure).
- Storage/retrieval facilities for elementary and parsed trees.
- The production of postscript files corresponding to elementary and parsed trees.
- Graphical displays of tree and feature data structures, including a scroll ‘web’ for large tree structures.
- Mouse-based tree editor for creating and modifying trees and feature structures.
- Hand combination of trees by adjunction or substitution for use in diagnosing grammar problems.
- Metarule tool for automatic aid to the generation of trees by using tree-based transformation rules

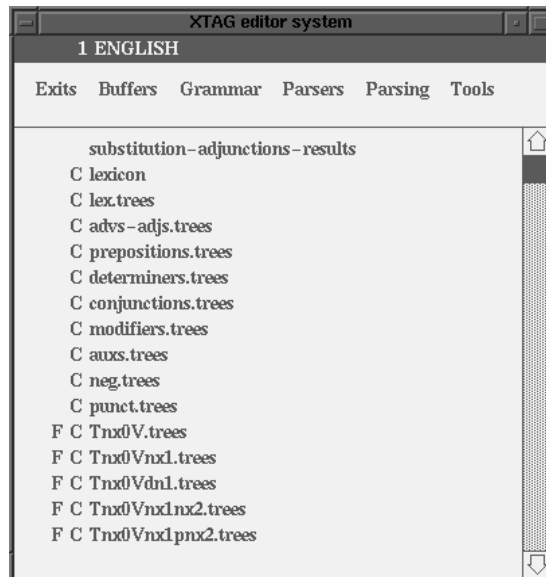


Figure 3.4: Interface to the XTAG system

3.2 Computer Platform

XTAG was developed on the Sun SPARC station series. It has been tested on various Sun platforms including Ultra-1, Ultra-Enterprise. XTAG is freely available from the XTAG web page at <http://www.cis.upenn.edu/~xtag/>. It requires 75 MB of disk space (once all binaries and databases are created after the install). XTAG requires the following software to run:

- A machine running UNIX and X11R4 (or higher). Previous releases of X will not work. X11R4 is free software which usually comes bundled with your OS. It is also freely available for various platforms at <http://www.xfree86.org/>
- A Common Lisp compiler which supports the latest definition of Common Lisp (Steele's Common Lisp, second edition). XTAG has been tested on Lucid Common Lisp/SPARC Solaris, Version: 4.2.1. Allegro CL is no longer directly supported, however there have been third party ports to recent versions of Allegro CL.
- CLX version 4 or higher. CLX is the Lisp equivalent to the Xlib package written in C.
- Mark Kantrowitz's Lisp Utilities from CMU: logical-pathnames and defsystem.

A patched version of CLX (Version 5.02) for SunOS 5.5.1 and the CMU Lisp Utilities are provided in our ftp directory for your convenience. However, we ask that you refer to the appropriate sources for updates.

The morphology database component ([Karp *et al.*, 1992]), no longer under licensing restrictions, is available as a separate download from the XTAG web page (see above for URL).

The syntactic database component is also available as part of the XTAG system ([Egedi and Martin, 1994]).

CHAPTER 3. OVERVIEW OF THE XTAG SYSTEM

More information can be obtained on the XTAG web page at
<http://www.cis.upenn.edu/~xtag/>.

Chapter 4

Underview

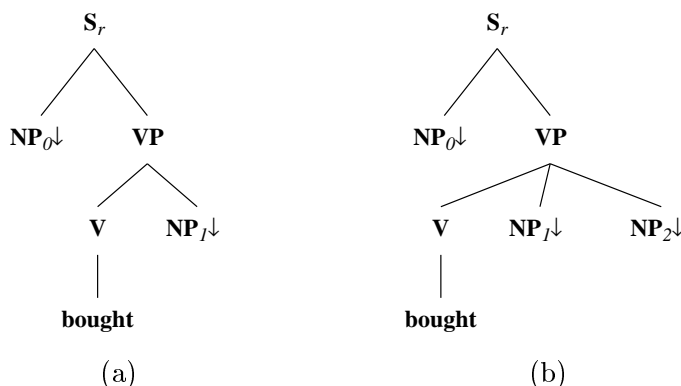
The morphology, syntactic, and tree databases together comprise the English grammar. A lexical item that is not in the databases receives a default tree selection and features for its part of speech and morphology. In designing the grammar, a decision was made early on to err on the side of acceptance whenever there are conflicting opinions as to whether or not a construction is grammatical. In this sense, the XTAG English grammar is intended to function primarily as an acceptor rather than a generator of English sentences. The range of syntactic phenomena that can be handled is large and includes auxiliaries (including inversion), copula, raising and small clause constructions, topicalization, relative clauses, infinitives, gerunds, passives, adjuncts, it-clefts, wh-clefts, PRO constructions, noun-noun modifications, extraposition, determiner sequences, genitives, negation, noun-verb contractions, clausal adjuncts and imperatives.

4.1 Subcategorization Frames

Elementary trees for non-auxiliary verbs are used to represent the linguistic notion of subcategorization frames. The anchor of the elementary tree subcategorizes for the other elements that appear in the tree, forming a clausal or sentential structure. Tree families group together trees belonging to the same subcategorization frame. Consider the following uses of the verb *buy*:

- (1) Srinu bought a book.
- (2) Srinu bought Beth a book.

In sentence (1), the verb *buy* subcategorizes for a direct object NP. The elementary tree anchored by *buy* is shown in Figure 4.1(a) and includes nodes for the NP complement of *buy* and for the NP subject. In addition to this declarative tree structure, the tree family also contains the trees that would be related to each other transformationally in a movement based approach, i.e. passivization, imperatives, wh-questions, relative clauses, and so forth. Sentence (2) shows that *buy* also subcategorizes for a double NP object. This means that *buy* also selects the double NP object subcategorization frame, or tree family, with its own set of transformationally related sentence structures. Figure 4.1(b) shows the declarative structure for this set of sentence structures.

Figure 4.1: Different subcategorization frames for the verb *buy*

4.2 Complements and Adjuncts

Complements and adjuncts have very different structures in the XTAG grammar. Complements are included in the elementary tree anchored by the verb that selects them, while adjuncts do not originate in the same elementary tree as the verb anchoring the sentence, but are instead added to a structure by adjunction. The contrasts between complements and adjuncts have been extensively discussed in the linguistics literature and the classification of a given element as one or the other remains a matter of debate (see [Rizzi, 1990], [Larson, 1988], [Jackendoff, 1990], [Larson, 1990], [Cinque, 1990], [Obernauer, 1984], [Lasnik and Saito, 1984], and [Chomsky, 1986]). The guiding rule used in developing the XTAG grammar is whether or not the sentence is ungrammatical without the questioned structure.¹ Consider the following sentences:

- (3) Srimi bought a book.
- (4) Srimi bought a book at the bookstore.
- (5) Srimi arranged for a ride.
- (6) *Srimi arranged.

Prepositional phrases frequently occur as adjuncts, and when they are used as adjuncts they have a tree structure such as that shown in Figure 4.2(a). This adjunction tree would adjoin into the tree shown in Figure 4.1(a) to generate sentence (4). There are verbs, however, such as *arrange*, *hunger* and *differentiate*, that take prepositional phrases as complements. Sentences (5) and (6) clearly show that the prepositional phrase are not optional for *arrange*. For these sentences, the prepositional phrase will be an initial tree (as shown in Figure 4.2(b)) that substitutes into an elementary tree, such as the one anchored by the verb *arrange* in Figure 4.2(c).

Virtually all parts of speech, except for main verbs, function as both complements and adjuncts in the grammar. More information is available in this report on various parts of

¹Iteration of a structure can also be used as a diagnostic: *Srimi bought a book at the bookstore on Walnut Street for a friend*.

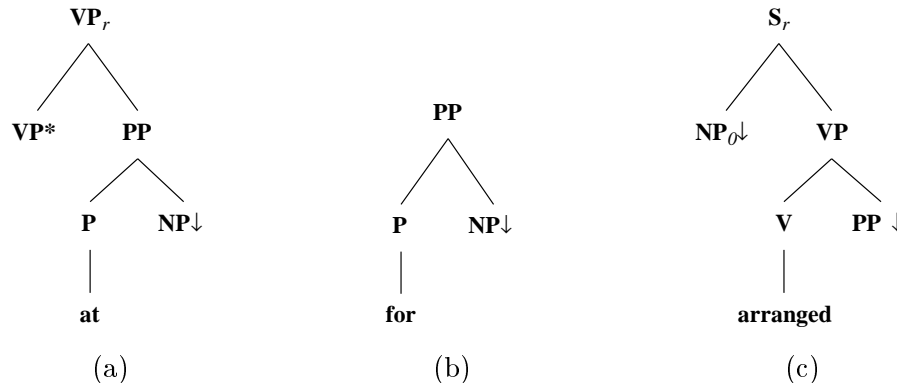


Figure 4.2: Trees illustrating the difference between Complements and Adjuncts

speech as complements: adjectives (e.g. section 6.13), nouns (e.g. section 6.2), and prepositions (e.g. section 6.10); and as adjuncts: adjectives (section 19.1), adverbs (section 19.5), nouns (section 19.2), and prepositions (section 19.4).

4.3 Non-S constituents

Although sentential trees are generally considered to be special cases in any grammar, insofar as they make up a ‘starting category’, it is the case that any initial tree constitutes a phrasal constituent. These initial trees may have substitution nodes that need to be filled (by other initial trees), and may be modified by adjunct trees, exactly as the trees rooted in S. Although grouping is possible according to the heads or anchors of these trees, we have not found any classification similar to the subcategorization frames for verbs that can be used by a lexical entry to ‘group select’ a set of trees. These trees are selected one by one by each lexical item, according to each lexical item’s idiosyncrasies. The grammar described by this technical report places them into several files for ease of use, but these files do not constitute tree families in the way that the subcategorization frames do.

4.4 Case Assignment

4.4.1 Approaches to Case

4.4.1.1 Case in GB theory

GB (Government and Binding) theory proposes the following ‘case filter’ as a requirement on S-structure.²

CASE FILTER Every overt NP must be assigned abstract case. [Haegeman, 1991]

²There are certain problems with applying the case filter as a requirement at the level of S-structure. These issues are not crucial to the discussion of the English XTAG implementation of case and so will not be discussed here. Interested readers are referred to [Lasnik and Uriagereka, 1988].

Abstract case is taken to be universal. Languages with rich morphological case marking, such as Latin, and languages with very limited morphological case marking, like English, are all presumed to have full systems of abstract case that differ only in the extent of morphological realization.

In GB, abstract case is argued to be assigned to NP's by various case assigners, namely verbs, prepositions, and INFL. Verbs and prepositions are said to assign accusative case to NP's that they govern, and INFL assigns nominative case to NP's that it governs. These governing categories are constrained as to where they can assign case by means of 'barriers' based on 'minimality conditions', although these are relaxed in 'exceptional case marking' situations. The details of the GB analysis are beyond the scope of this technical report, but see [Chomsky, 1986] for the original analysis or [Haegeman, 1991] for an overview. Let it suffice for us to say that the notion of abstract case and the case filter are useful in accounting for a number of phenomena including the distribution of nominative and accusative case, and the distribution of overt NP's and empty categories (such as PRO).

4.4.1.2 Minimalism and Case

A major conceptual difference between GB theories and Minimalism is that in Minimalism, lexical items carry their features with them rather than being assigned their features based on the nodes that they end up at. For nouns, this means that they carry case with them, and that their case is 'checked' when they are in SPEC position of AGR_s or AGR_o , which subsequently disappears [Chomsky, 1992].

4.4.2 Case in XTAG

The English XTAG grammar adopts the notion of case and the case filter for many of the same reasons argued in the GB literature. However, in some respects the English XTAG grammar's implementation of case more closely resembles the treatment in Chomsky's Minimalism framework [Chomsky, 1992] than the system outlined in the GB literature [Chomsky, 1986]. As in Minimalism, nouns in the XTAG grammar carry case with them, which is eventually 'checked'. However in the XTAG grammar, noun cases are checked against the case values assigned by the verb during the unification of the feature structures. Unlike Chomsky's Minimalism, there are no separate AGR nodes; the case checking comes from the verbs directly. Case assignment from the verb is more like the GB approach than the requirement of a SPEC-head relationship in Minimalism.

Most nouns in English do not have separate forms for nominative and accusative case, and so they are ambiguous between the two. Pronouns, of course, are morphologically marked for case, and each carries the appropriate case in its feature. Figures 4.3(a) and 4.3(b) show the NP tree anchored by a noun and a pronoun, respectively, along with the feature values associated with each word. Note that *books* simply gets the default case **nom/acc**, while *she* restricts the case to be **nom**.

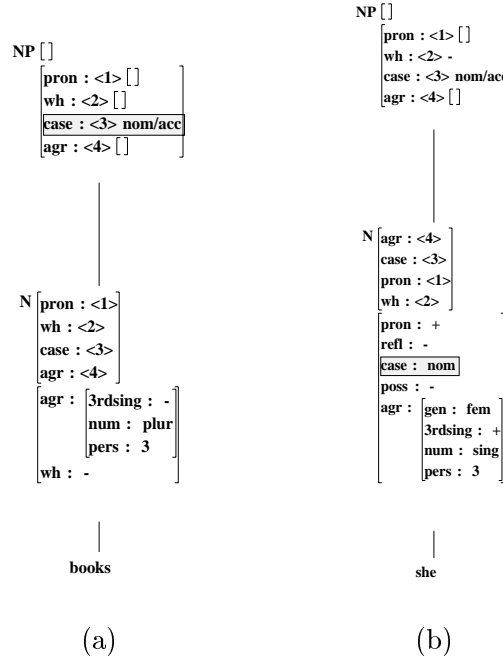


Figure 4.3: Lexicalized NP trees with case markings

4.4.3 Case Assigners

4.4.3.1 Prepositions

Case is assigned in the XTAG English grammar by two lexical categories - verbs and prepositions.³ Prepositions assign accusative case (**acc**) through their **<assign-case>** feature, which is linked directly to the **<case>** feature of their objects. Figure 4.4(a) shows a lexicalized preposition tree, while Figure 4.4(b) shows the same tree with the NP tree from Figure 4.3(a) substituted into the NP position. Figure 4.4(c) is the tree in Figure 4.4(b) after unification has taken place. Note that the case ambiguity of *books* has been resolved to accusative case.

4.4.3.2 Verbs

Verbs are the other part of speech in the XTAG grammar that can assign case. Because XTAG does not distinguish INFL and VP nodes, verbs must provide case assignment on the subject position in addition to the case assigned to their NP complements.

Assigning case to NP complements is handled by building the case values of the complements directly into the tree that the case assigner (the verb) anchors. Figures 4.5(a) and 4.5(b) show an S tree⁴ that would be anchored⁵ by a transitive and ditransitive verb, respectively. Note that the case assignments for the NP complements are already in the tree, even though there is not yet a lexical item anchoring the tree. Since every verb that selects these trees (and other

³For also assigns case as a complementizer. See section 8.5 for more details.

⁴Features not pertaining to this discussion have been taken out to improve readability and to make the trees easier to fit onto the page.

⁵The diamond marker (\diamond) indicates the anchor(s) of a structure if the tree has not yet been lexicalized.

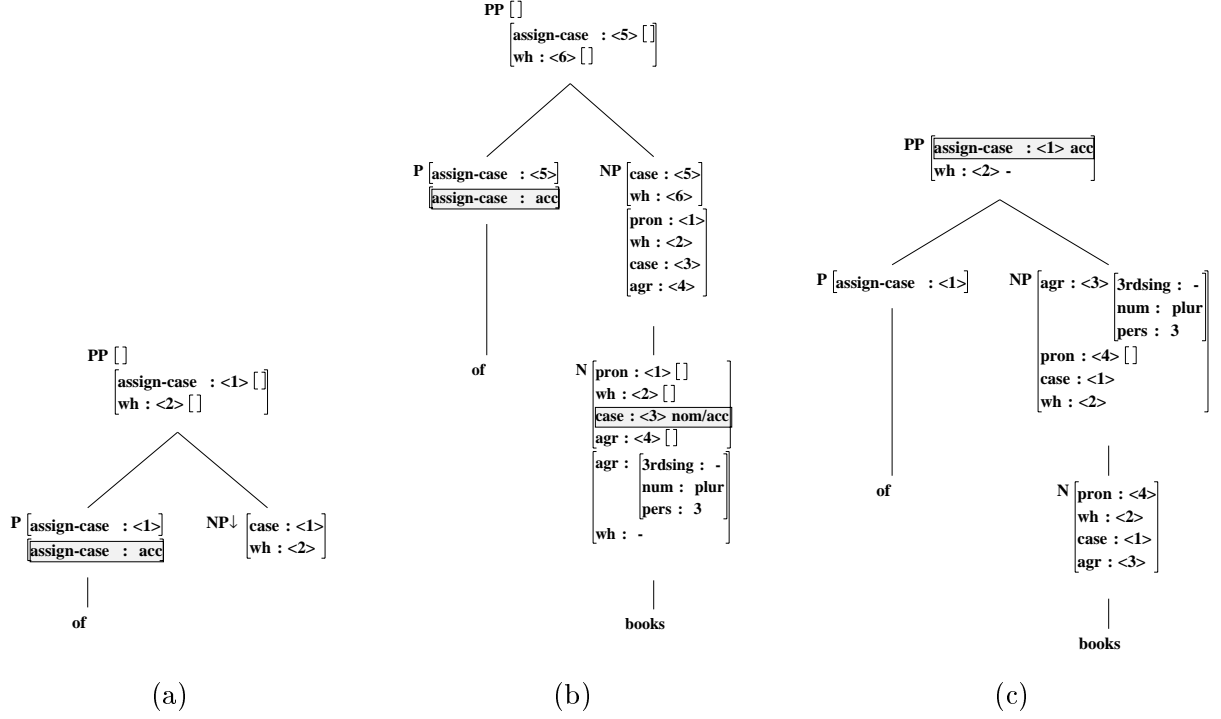


Figure 4.4: Assigning case in prepositional phrases

trees in each respective subcategorization frame) assigns the same case to the complements, building case features into the tree has exactly the same result as putting the case feature value in each verb's lexical entry.

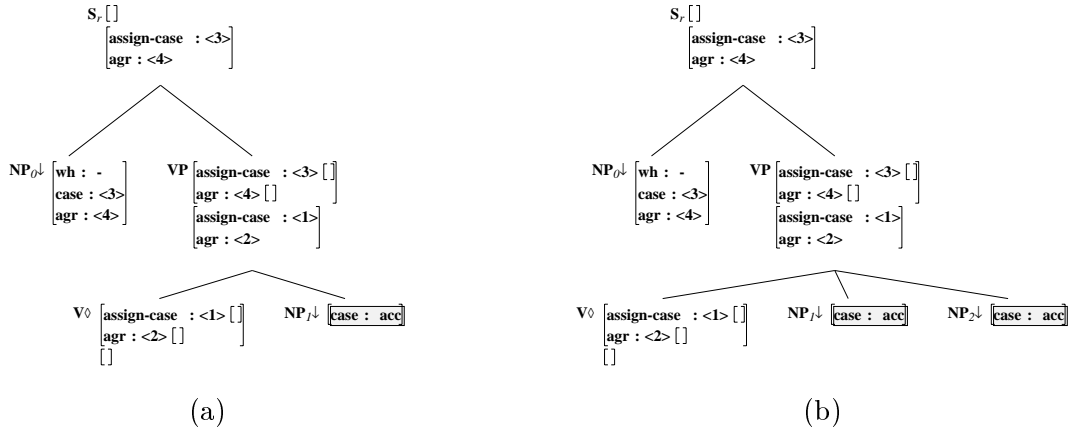


Figure 4.5: Case assignment to NP arguments

The case assigned to the subject position varies with verb form. Since the XTAG grammar treats the inflected verb as a single unit rather than dividing it into INFL and V nodes, case, along with tense and agreement, is expressed in the features of verbs, and must be passed in the appropriate manner. The trees in Figure 4.6 show the path of linkages that joins the **<assign-**

case> feature of the V to the <**case**> feature of the subject NP. The morphological form of the verb determines the value of the <**assign-case**> feature. Figures 4.6(a) and 4.6(b) show the same tree⁶ anchored by different morphological forms of the verb *sing*, which give different values for the <**assign-case**> feature.

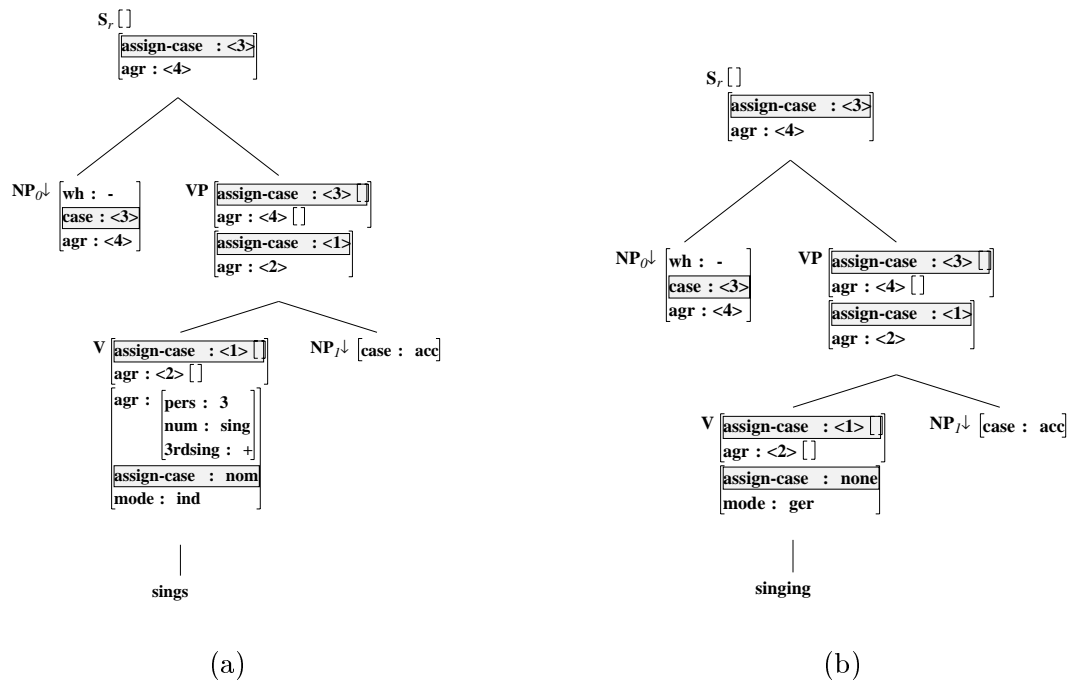


Figure 4.6: Assigning case according to verb form

The adjunction of an auxiliary verb onto the VP node breaks the <**assign-case**> link from the main V, replacing it with a link from the auxiliary verb instead.⁷ The progressive form of the verb in Figure 4.6(b) has the feature-value <**assign-case**>=**none**, but this is overridden by the adjunction of the appropriate form of the auxiliary word *be*. Figure 4.7(a) shows the lexicalized auxiliary tree, while Figure 4.7(b) shows it adjoined into the transitive tree shown in Figure 4.6(b). The case value passed to the subject NP is now **nom** (nominative).

4.4.4 PRO in a unification based framework

Tensed forms of a verb assign nominative case, and untensed forms assign case **none**, as the progressive form of the verb *sing* does in Figure 4.6(b). This is different than assigning no case at all, as one form of the infinitive marker *to* does. See Section 8.5 for more discussion of this special case.) The distinction of a case **none** from no case is indicative of a divergence from the standard GB theory. In GB theory, the absence of case on an NP means that only PRO can fill that NP. With feature unification as is used in the FB-LTAG grammar, the absence of case on an NP means that *any* NP can fill it, regardless of its case. This is due to the mechanism of unification, in which if something is unspecified, it can unify with anything. Thus we have

⁶ Again, the feature structures shown have been restricted to those that pertain to the V/NP interaction.

⁷ See section 20.1 for a more complete explanation of how this relinking occurs.

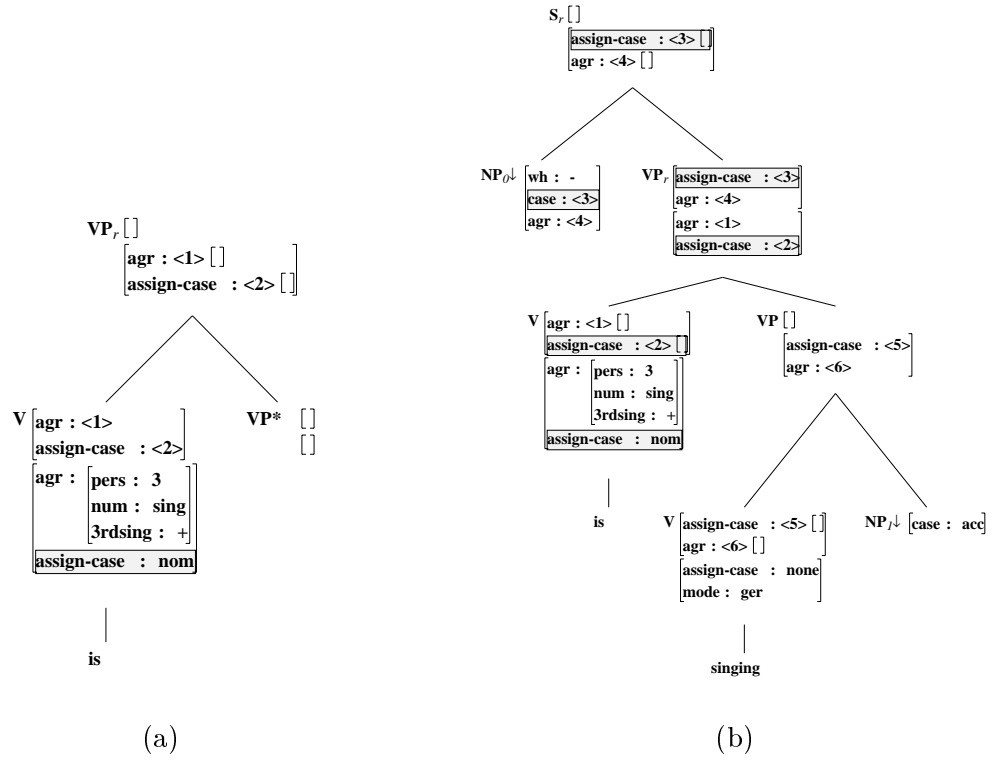


Figure 4.7: Proper case assignment with auxiliary verbs

a specific case **none** to handle verb forms that in GB theory do not assign case. PRO is the only NP with case **none**. Note that although we are drawn to this treatment by our use of unification for feature manipulation, our treatment is very similar to the assignment of null case to PRO in [Chomsky and Lasnik, 1993]. [Watanabe, 1993] also proposes a very similar approach within Chomsky's Minimalist framework.⁸

⁸See Sections 8.1 and 8.9 for additional discussion of PRO.

CHAPTER 4. UNDERVIEW

Part II

Verb Classes

Chapter 5

Where to Find What

The two page table that follows gives an overview of what types of trees occur in various tree families with pointers to discussion in this report. An entry in a cell of the table indicates that the tree(s) for the construction named in the row header are included in the tree family named in the column header. Entries are of two types. If the particular tree(s) are displayed and/or discussed in this report the entry gives a page number reference to the relevant discussion or figure.¹ Otherwise, a \checkmark indicates inclusion in the tree family but no figure or discussion related specifically to that tree in this report. Blank cells indicate that there are no trees for the construction named in the row header in the tree family named in the column header. Two tables are given below. The first one gives the expansion of abbreviations in the table headers. The second table gives the name given to each tree family in the actual XTAG grammar. This makes it easier to find the description of each tree family in Chapter 6 and to compare the description with the online XTAG grammar.

¹Since Chapter 6 has a brief discussion and a declarative tree for every tree family, page references are given only for other sections in which discussion or tree diagrams appear.

CHAPTER 5. WHERE TO FIND WHAT

Abbreviation	Full Name
Sent. Subj. w. <i>to</i>	Sentential Subject with <i>to</i> PP complement
Pred. Mult-wd. ARB, P	Predicative Multi-word PP with Adv, Prep anchors
Pred. Mult-wd. A, P	Predicative Multi-word PP with Adj, Prep anchors
Pred. Mult-wd. N, P	Predicative Multi-word PP with Noun, Prep anchors
Pred. Mult-wd. P, P	Predicative Multi-word PP with two Prep anchors
Pred. Mult-wd. no int. mod.	Predicative Multi-word PP with no internal modification
Pred. Sent. Subj., ARB, P	Predicative PP with Sentential Subject, and Adv, Prep anchors
Pred. Sent. Subj., A, P	Predicative PP with Sentential Subject, and Adj, Prep anchors
Pred. Sent. Subj., Conj, P	Predicative PP with Sentential Subject, and Conj, Prep anchors
Pred. Sent. Subj., N, P	Predicative PP with Sentential Subject, and Noun, Prep anchors
Pred. Sent. Subj., P, P	Predicative PP with Sentential Subject, and two Prep anchors
Pred. Sent. Subj., no int-mod	Predicative PP with Sentential Subject, no internal modification
Pred. Locative	Predicative anchored by a Locative Adverb
Pred. A Sent. Subj., Comp.	Predicative Adjective with Sentential Subject and Complement
Sentential Comp. with NP	Sentential Complement with NP
Pred. Mult wd. V, P	Predicative Multi-word with Verb, Prep anchors
Adj. Sm. Cl. w. Sentential Subj.	Adjective Small Clause with Sentential Subject
NP Sm. Clause w. Sentential Subj.	NP Small Clause with Sentential Subject
PP Sm. Clause w. Sentential Subj.	PP Small Clause with Sentential Subject
NP Sm. Cl. w. Sent. Comp.	NP Small Clause with Sentential Complement
Adj. Sm. Cl. w. Sent. Comp.	Adjective Small Clause with Sentential Complement
Exhaustive PP Sm. Cl.	Exhaustive PP Small Clause
Ditrans. Light Verbs w. PP Shift	Ditransitive Light Verbs with PP Shift
Ditrans. Light Verbs w/o PP Shift	Ditransitive Light Verbs without PP Shift
Y/N question	Yes/No question
Wh-mov. NP complement	Wh-moved NP complement
Wh-mov. S comp.	Wh-moved S complement
Wh-mov. Adj comp.	Wh-moved Adjective complement
Wh-mov. object of a P	Wh-moved object of a P
Wh-mov. PP	Wh-moved PP
Topic. NP complement	Topicalized NP complement
Det. gerund	Determiner gerund
Rel. cl. on NP comp.	Relative clause on NP complement
Rel. cl. on PP comp.	Relative clause on PP complement
Rel. cl. on NP object of P	Relative clause on NP object of P
Pass. with wh-moved subj.	Passive with wh-moved subject (with and without <i>by</i> phrase)
Pass. w. wh-mov. ind. obj.	Passive with wh-moved indirect object (with and without <i>by</i> phrase)
Pass. w. wh-mov. obj. of the <i>by</i> phrase	Passive with wh-moved object of the <i>by</i> phrase
Pass. w. wh-mov. <i>by</i> phrase	Passive with wh-moved <i>by</i> phrase
Trans. Idiom with V, D and N	Transitive Idiom with Verb, Det and Noun anchors
Idiom with V, D, N	Idiom with V, D, and N anchors
Idiom with V, D, A, N	Idiom with V, D, A, and N anchors
Idiom with V, N	Idiom with V, and N anchor
Idiom with V, A, N	Idiom with V, A, and N anchors
Idiom with V, D, N, P	Idiom with V, D, N, and Prep anchors
Idiom with V, D, A, N, P	Idiom with V, D, A, N, and Prep anchors
Idiom with V, N, P	Idiom with V, N, and Prep anchors
Idiom with V, A, N, P	Idiom with V, A, N, and Prep anchors

Full Name	XTAG Name
Intransitive Sentential Subject	Ts0V
Sentential Subject with ‘to’ complement	Ts0Vtonx1
PP Small Clause, with Adv and Prep anchors	Tnx0ARBPNx1
PP Small Clause, with Adj and Prep anchors	Tnx0APnx1
PP Small Clause, with Noun and Prep anchors	Tnx0NPNx1
PP Small Clause, with Prep anchors	Tnx0PPnx1
PP Small Clause, with Prep and Noun anchors	Tnx0PNaPNx1
PP Small Clause with Sentential Subject, and Adv and Prep anchors	Ts0ARBPNx1
PP Small Clause with Sentential Subject, and Adj and Prep anchors	Ts0APnx1
PP Small Clause with Sentential Subject, and Noun and Prep anchors	Ts0NPNx1
PP Small Clause with Sentential Subject, and Prep anchors	Ts0PPnx1
PP Small Clause with Sentential Subject, and Prep and Noun anchors	Ts0PNaPNx1
Exceptional Case Marking	TXnx0Vs1
Locative Small Clause with Ad anchor	Tnx0nx1ARB
Predicative Adjective with Sentential Subject and Complement	Ts0A1s1
Transitive	Tnx0Vnx1
Ditransitive with PP shift	Tnx0Vnx1tonx2
Ditransitive	Tnx0Vnx1nx2
Ditransitive with PP	Tnx0Vnx1pnx2
Sentential Complement with NP	Tnx0Vnx1s2
Intransitive Verb Particle	Tnx0Vpl
Transitive Verb Particle	Tnx0Vplnx1
Ditransitive Verb Particle	Tnx0Vplnx1nx2
Intransitive with PP	Tnx0Vpnx1
Sentential Complement	Tnx0Vs1
Light Verbs	Tnx0lVN1
Ditransitive Light Verbs with PP Shift	Tnx0lVN1Pnx2
Adjective Small Clause with Sentential Subject	Ts0Ax1
NP Small Clause with Sentential Subject	Ts0N1
PP Small Clause with Sentential Subject	Ts0Pnx1
Predicative Multi-word with Verb, Prep anchors	Tnx0VPnx1
Adverb It-Cleft	TIItVad1s2
NP It-Cleft	TIItVnx1s2
PP It-Cleft	TIItVpnx1s2
Adjective Small Clause Tree	Tnx0Ax1
Adjective Small Clause with Sentential Complement	Tnx0A1s1
Equative <i>BE</i>	Tnx0BEnx1
NP Small Clause	Tnx0N1
NP with Sentential Complement Small Clause	Tnx0N1s1
PP Small Clause	Tnx0Pnx1
Exhaustive PP Small Clause	Tnx0Px1
Intransitive	Tnx0V
Intransitive with Adjective	Tnx0Vax1
Transitive Sentential Subject	Ts0Vnx1
Idiom with V, D and N	Tnx0VDN1
Idiom with V, D, A, and N anchors	Tnx0VDAN1
Idiom with V and N anchors	Tnx0VN1
Idiom with V, A, and N anchors	Tnx0VAN1
Idiom with V, D, N, and Prep anchors	Tnx0VDN1Pnx2
Idiom with V, D, A, N, and Prep anchors	Tnx0VDAN1Pnx2
Idiom with V, N, and Prep anchors	Tnx0VN1Pnx2
Idiom with V, A, N, and Prep anchors	Tnx0VAN1Pnx2

Constructions	Tree families														
	Intransitive Sentential Subj	Sent. Subj. w. to	Pred. Mult-wd. ARB, P	Pred. Mult-wd. A, P	Pred. Mult-wd. N, P	Pred. Mult-wd. P, P	Pred. Mult-wd. no int. mod.	Pred. Sent. Subj., ARB, P	Pred. Sent. Subj., A, P	Pred. Sent. Subj., N, P	Pred. Sent. Subj., P, P	Pred. Sent. Subj., no int-mod	ECM	Pred. Locative	Pred. A Sent. Subj., Comp.
Declarative	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	89	71	
Passive w/ & w/o <i>by</i> phrase													90		
Y/N quest.															
Wh-moved subject	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Wh-mov. NP complement, DO or IO															
Wh-mov. S comp.															
Wh-mov. Adj. or Adv. comp.														72	
Wh-mov. object of a P			✓	✓	✓	✓	✓								
Wh-mov. PP			✓			✓	✓								
Topic. NP comp.															
Imperative			✓			✓	✓						✓	✓	
Det. gerund															
NP gerund			✓				✓						✓	✓	
Ergative															
Rel. cl. on subj. w/ NP			✓	✓	✓	✓	✓						✓	✓	
Rel. cl. on subj. w/ Comp			✓	✓	✓	✓	✓						✓	✓	
Rel. cl. on NP comp., DO, IO w/ NP															
Rel. cl. on NP comp., DO, IO w/ Comp															
Rel. cl. on PP comp. w/ pied-piping						✓	✓								
Rel. cl. on NP object of P w/ NP			✓	✓	✓	✓	✓	✓	✓	✓	✓	✓			
Rel. cl. on NP object of P w/ Comp			✓	✓		✓	✓	✓	✓	✓	✓	✓			
Rel. cl. on adjunct w/ PP	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓		✓
Rel. cl. on adjunct w/ Comp	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓		✓
Pass. w. wh-mov. subj.															
Pass. w. wh-mov. ind. obj.															
Pass. w. wh-mov. obj. of <i>by</i> phrase															
Pass. w. wh-mov. <i>by</i> phrase															

Constructions	Tree families															
	Transitive	Ditransitive with PP shift	Ditransitive	Ditransitive with PP	Sentential Comp. with NP	Intransitive Verb Particle	Transitive Verb Particle	Ditransitive Verb Particle	Intransitive with PP	Sentential Complement	Trans. Light Vs	Ditrans. Light Vs	Adj. Sm. Cl. w. Sentential Subj.	NP Sm. Cl. w. Sentential Subj.	PP Sm. Cl. w. Sentential Subj.	Pred. Mult. wd. V, P
Declarative	24	108	24	✓	✓	✓	✓	✓	21	10,86	✓	✓	✓	✓	✓	✓
Passive w/ & w/o <i>by</i> phrase	✓	✓	✓	✓	115		✓	✓				✓				✓
Y/N quest.																
Wh-moved subject	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓		✓	✓	✓	✓
Wh-mov. NP complement, DO or IO	118	✓	122	✓	✓		✓	✓								
Wh-mov. S comp.					✓					✓						
Wh-mov. Adj. or Adv. comp.													✓			
Wh-mov. object of a P		✓		122					✓			✓				✓
Wh-mov. PP		✓		123					✓			✓				
Topic. NP comp.	✓	✓	✓	✓	✓		✓	✓								
Imperative	141	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓				✓
Det. gerund	144	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓				✓
NP gerund	145	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓				✓
Ergative	80															
Rel. cl. on subj. w/ NP	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓				✓
Rel. cl. on subj. w/ Comp	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓				✓
Rel. cl. on NP comp., DO, IO w/ NP	✓	✓	✓	✓	✓		✓	✓	✓			✓				
Rel. cl. on NP comp., DO, IO w/ Comp	✓	✓	✓	✓	✓		✓	✓	✓			✓				
Rel. cl. on PP comp. w/ pied-piping	✓	✓	✓	✓	✓		✓		✓			✓			✓	
Rel. cl. on NP object of P w/ NP	✓	✓	✓	✓	✓		✓		✓			✓			✓	✓
Rel. cl. on NP object of P w/ Comp	✓	✓	✓	✓	✓		✓		✓			✓			✓	✓
Rel. cl. on adjunct w/ PP	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Rel. cl. on adjunct w/ Comp	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Parenthetical quoting clause					✓					✓						
Past-participial as arg Adj	✓															✓
Past-participial NP pre-mod	✓															✓
Pass. w. wh-mov. subj.	✓	✓	✓	✓	✓		✓	✓				✓				✓
Pass. w. wh-mov. ind. obj.		✓	✓	✓	✓			✓				✓				
Pass. w. wh-mov. obj. of <i>by</i> phrase	✓	✓	✓	✓	✓		✓	✓				✓				✓
Pass. w. wh-mov. <i>by</i> phrase	✓	✓	✓	✓	✓		✓	✓								✓

Constructions	Tree families												
	Adverb It-Cleft	NP It-Cleft	PP It-Cleft	Adj. Small Clause	Adj. Sm. Cl. w. Sent. Comp.	Equative <i>BE</i>	NP Small Clause	NP Sm. Cl. w. Sent. Comp.	PP Small Clause	Exhaustive PP Sm. Cl.	Intransitive	Intransitive with Adjective	Transitive Sentential Subj
Declarative	✓	✓	111	104	✓	106	101	✓	101	✓	✓	✓	91
Passive w/ & w/o <i>by</i> phrase													
Y/N quest.	✓	✓	111			✓							
Wh-moved subject				87	✓		✓	✓	✓	✓	121	87	✓
Wh-mov. NP complement, DO or IO		✓					✓						✓
Wh-mov. S comp.													
Wh-mov. Adj. or Adv. comp.	✓			✓								124	
Wh-mov. object of a P									✓				
Wh-mov. PP			✓						✓				
Topic. NP comp.		✓					✓						
Imperative				✓	✓		✓	✓	✓	✓	✓	✓	
Det. gerund													
NP gerund				✓	✓		✓	✓	✓	✓	✓	✓	
Ergative													
Rel. cl. on subj. w/ NP				✓	✓		✓	✓	✓	✓	✓	✓	
Rel. cl. on subj. w/ Comp				✓	✓		✓	✓	✓	✓	✓	✓	
Rel. cl. on NP comp., DO, IO w/ NP													
Rel. cl. on NP comp., DO, IO w/ Comp													
Rel. cl. on PP comp. w/ pied-piping									✓				
Rel. cl. on NP object of P w/ NP									✓				
Rel. cl. on NP object of P w/ Comp									✓				
Rel. cl. on adjunct w/ PP	✓	✓	✓	✓	✓		✓	✓	✓	✓	✓	✓	✓
Rel. cl. on adjunct w/ Comp	✓	✓	✓	✓	✓		✓	✓	✓	✓	✓	✓	✓
Participial NP pre-mod											✓		
Pass. w. wh-mov. subj.													
Pass. w. wh-mov. ind. obj.													
Pass. w. wh-mov. obj. of <i>by</i> phrase													
Pass. w. wh-mov. <i>by</i> phrase													

Constructions	Tree families							
	Idiom with V, D, N	Idiom with V, D, A, N	Idiom with V, N	Idiom with V, A, N	Idiom with V, D, N, P	Idiom with V, D, A, N, P	Idiom with V, N, P	Idiom with V, A, N, P
Declarative	✓	✓	✓	✓	✓	✓	✓	✓
Passive w/ & w/o <i>by</i> phrase	✓	✓	✓	✓	✓	✓	✓	✓
Y/N quest.								
Wh-moved subject	✓	✓	✓	✓	✓	✓	✓	✓
Wh-mov. NP complement, DO or IO								
Wh-mov. S comp.								
Wh-mov. Adj. or Adv. comp.								
Wh-mov. object of a P								
Wh-mov. PP								
Topic. NP comp.								
Imperative	✓	✓	✓	✓	✓	✓	✓	✓
Det. gerund								
NP gerund	✓	✓	✓	✓	✓	✓	✓	✓
Ergative								
Rel. cl. on subj. w/ NP	✓	✓	✓	✓	✓	✓	✓	✓
Rel. cl. on subj. w/ Comp	✓	✓	✓	✓	✓	✓	✓	✓
Rel. cl. on NP comp., DO, IO w/ NP								
Rel. cl. on NP comp., DO, IO w/ Comp								
Rel. cl. on PP comp. w/ pied-piping								
Rel. cl. on NP object of P w/ NP								
Rel. cl. on NP object of P w/ Comp								
Rel. cl. on adjunct w/ PP	✓	✓	✓	✓	✓	✓	✓	✓
Rel. cl. on adjunct w/ Comp	✓	✓	✓	✓	✓	✓	✓	✓
Pass. w. wh-mov. subj.								
Pass. w. wh-mov. ind. obj.								
Pass. w. wh-mov. obj. of <i>by</i> phrase	✓	✓	✓	✓	✓	✓	✓	✓
Pass. w. wh-mov. <i>by</i> phrase	✓	✓	✓	✓	✓	✓	✓	✓
Outer Pass. w. and wo. <i>by</i> phrase					✓	✓	✓	✓
Outer Pass. w. Rel. cl. on subj. w. Comp	✓				✓	✓	✓	✓
Outer Pass. w. Rel. cl. on subj. w. NP	✓				✓	✓	✓	✓

Chapter 6

Verb Classes

Each main¹ verb in the syntactic lexicon selects at least one tree family² (subcategorization frame). Since the tree database and syntactic lexicon are already separated for space efficiency (see Chapter 3), each verb can efficiently select a large number of trees by specifying a tree family, as opposed to each of the individual trees. This approach allows for a considerable reduction in the number of trees that must be specified for any given verb or form of a verb.

There are currently 52 tree families in the system.³ This chapter gives a brief description of each tree family and shows the corresponding declarative tree⁴, along with any peculiar characteristics or trees. It also indicates which transformations are in each tree family, and gives the number of verbs that select that family.⁵ A few sample verbs are given, along with example sentences.

6.1 Intransitive: Tnx0V

Description: This tree family is selected by verbs that do not require an object complement of any type. Adverbs, prepositional phrases and other adjuncts may adjoin on, but are not required for the sentences to be grammatical. 1,878 verbs select this family.

Examples: *eat, sleep, dance*

Al ate .

Seth slept .

Hyun danced .

Declarative tree: See Figure 6.1.

Other available trees: wh-moved subject, subject relative clause with and without comp, adjunct (gap-less) relative clause with comp, adjunct (gap-less) relative clause with PP pied-piping, imperative, determiner gerund, NP gerund, pre-nominal participial.

¹Auxiliary verbs are handled under a different mechanism. See Chapter 20 for details.

²See section 3.1.2 for explanation of tree families.

³An explanation of the naming convention used in naming the trees and tree families is available in Appendix D.

⁴Before lexicalization, the \diamond indicates the anchor of the tree.

⁵Numbers given are as of August 1998 and are subject to some change with further development of the grammar.

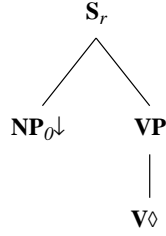


Figure 6.1: Declarative Intransitive Tree: $\alpha nx0V$

6.2 Transitive: Tnx0Vnx1

Description: This tree family is selected by verbs that require only an NP object complement. The NP's may be complex structures, including gerund NP's and NP's that take sentential complements. This does not include light verb constructions (see sections 6.15 and 6.16). 4,343 verbs select the transitive tree family.

Examples: *eat, dance, take, like*

Al ate an apple .

Seth danced the tango .

Hyun is taking an algorithms course .

Anoop likes the fact that the semester is finished .

Declarative tree: See Figure 6.2.

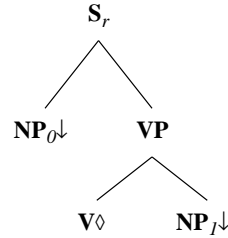


Figure 6.2: Declarative Transitive Tree: $\alpha nx0Vnx1$

Other available trees: wh-moved subject, wh-moved object, subject relative clause with and without comp, adjunct (gap-less) relative clause with comp/with PP pied-piping, object relative clause with and without comp, imperative, determiner gerund, NP gerund, passive with *by* phrase, passive without *by* phrase, passive with wh-moved subject and *by* phrase, passive with wh-moved subject and no *by* phrase, passive with wh-moved object out of the *by* phrase, passive with wh-moved *by* phrase, passive with relative clause on subject and *by* phrase with and without comp, passive with relative clause on subject and no *by* phrase with and without comp, passive with relative clause on object on the *by* phrase with and without comp/with PP pied-piping, gerund passive with *by* phrase, gerund passive without *by* phrase, ergative, ergative with wh-moved subject, ergative with subject relative clause with and without comp, ergative with adjunct (gap-less) relative clause with comp/with

PP pied-piping. In addition, two other trees that allow transitive verbs to function as adjectives (e.g. *the stopped truck*) are also in the family.

6.3 Ditransitive: Tnx0Vnx1nx2

Description: This tree family is selected by verbs that take exactly two NP complements. It does **not** include verbs that undergo the ditransitive verb shift (see section 6.5). The apparent ditransitive alternates involving verbs in this class and benefactive PP's (e.g. *John baked a cake for Mary*) are analyzed as transitives (see section 6.2) with a PP adjunct. Benefactives are taken to be adjunct PP's because they are optional (e.g. *John baked a cake* vs. *John baked a cake for Mary*). 122 verbs select the ditransitive tree family.

Examples: *ask, cook, win*

Christy asked Mike a question .

Doug cooked his father dinner .

Dania won her sister a stuffed animal .

Declarative tree: See Figure 6.3.

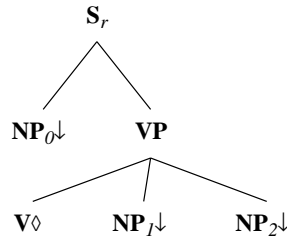


Figure 6.3: Declarative Ditransitive Tree: $\alpha nx0Vnx1nx2$

Other available trees: wh-moved subject, wh-moved direct object, wh-moved indirect object, subject relative clause with and without comp, adjunct (gap-less) relative clause with comp/with PP pied-piping, direct object relative clause with and without comp, indirect object relative clause with and without comp, imperative, determiner gerund, NP gerund, passive with *by* phrase, passive without *by* phrase, passive with wh-moved subject and *by* phrase, passive with wh-moved subject and no *by* phrase, passive with wh-moved object out of the *by* phrase, passive with wh-moved *by* phrase, passive with wh-moved indirect object and *by* phrase, passive with wh-moved indirect object and no *by* phrase, passive with relative clause on subject and *by* phrase with and without comp, passive with relative clause on subject and no *by* phrase with and without comp, passive with relative clause on object of the *by* phrase with and without comp/with PP pied-piping, passive with relative clause on the indirect object and *by* phrase with and without comp, passive with relative clause on the indirect object and no *by* phrase with and without comp, passive with/without *by*-phrase with adjunct (gap-less) relative clause with comp/with PP pied-piping, gerund passive with *by* phrase, gerund passive without *by* phrase.

6.4 Ditransitive with PP: Tnx0Vnx1pnx2

Description: This tree family is selected by ditransitive verbs that take a noun phrase followed by a prepositional phrase. The preposition is not constrained in the syntactic lexicon. The preposition must be required and not optional - that is, the sentence must be ungrammatical with just the noun phrase (e.g. **John put the table*). No verbs, therefore, should select both this tree family and the transitive tree family (see section 6.2). This tree family is also distinguished from the ditransitive verbs, such as *give*, that undergo verb shifting (see section 6.5). There are 62 verbs that select this tree family.

Examples: *associate, put, refer*
Rostenkowski associated money with power .
He put his reputation on the line .
He referred all questions to his attorney .

Declarative tree: See Figure 6.4.

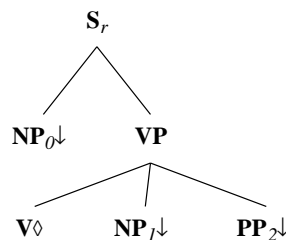


Figure 6.4: Declarative Ditransitive with PP Tree: $\alpha nx0Vnx1pnx2$

Other available trees: wh-moved subject, wh-moved direct object, wh-moved object of PP, wh-moved PP, subject relative clause with and without comp, adjunct (gap-less) relative clause with comp/with PP pied-piping, direct object relative clause with and without comp, object of PP relative clause with and without comp/with PP pied-piping, imperative, determiner gerund, NP gerund, passive with *by* phrase, passive without *by* phrase, passive with wh-moved subject and *by* phrase, passive with wh-moved subject and no *by* phrase, passive with wh-moved object out of the *by* phrase, passive with wh-moved *by* phrase, passive with wh-moved object out of the PP and *by* phrase, passive with wh-moved object out of the PP and no *by* phrase, passive with wh-moved PP and *by* phrase, passive with wh-moved PP and no *by* phrase, passive with relative clause on subject and *by* phrase with and without comp, passive with relative clause on subject and no *by* phrase with and without comp, passive with relative clause on object of the *by* phrase with and without comp/with PP pied-piping, passive with relative clause on the object of the PP and *by* phrase with and without comp/with PP pied-piping, passive with relative clause on the object of the PP and no *by* phrase with and without comp/with PP pied-piping, passive with and without *by* phrase with adjunct (gap-less) relative clause with comp/with PP pied-piping, gerund passive with *by* phrase, gerund passive without *by* phrase.

6.5 Ditransitive with PP shift: Tnx0Vnx1tonx2

Description: This tree family is selected by ditransitive verbs that undergo a shift to a *to* prepositional phrase. These ditransitive verbs are clearly constrained so that when they shift, the prepositional phrase must start with *to*. This is in contrast to the Ditransitives with PP in section 6.4, in which verbs may appear in [NP V NP PP] constructions with a variety of prepositions. Both the dative shifted and non-shifted PP complement trees are included. 56 verbs select this family.

Examples: *give, promise, tell*

Bill gave Hillary flowers .

Bill gave flowers to Hillary .

Whitman promised the voters a tax cut .

Whitman promised a tax cut to the voters .

Pinnocchio told Gepetto a lie .

Pinnocchio told a lie to Gepetto .

Declarative tree: See Figure 6.5.

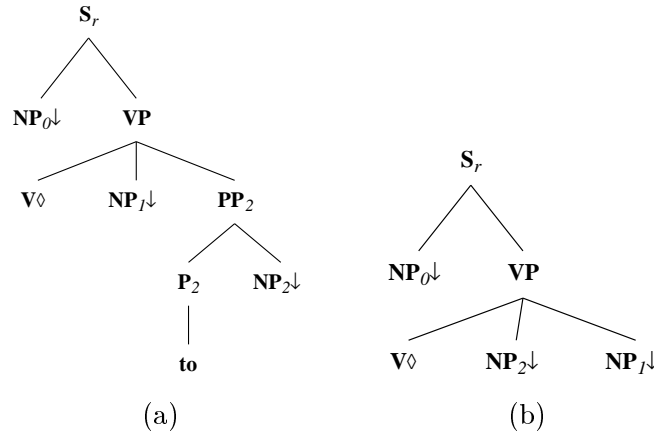


Figure 6.5: Declarative Ditransitive with PP shift Trees: $\alpha nx0Vnx1Pnx2$ (a) and $\alpha nx0Vnx2nx1$ (b)

Other available trees: Non-shifted: wh-moved subject, wh-moved direct object, wh-moved indirect object, subject relative clause with and without comp, adjunct (gap-less) relative clause with comp/with PP pied-piping, direct object relative clause with comp/with PP pied-piping, indirect object relative clause with and without comp/with PP pied-piping, imperative, NP gerund, passive with *by* phrase, passive without *by* phrase, passive with wh-moved subject and *by* phrase, passive with wh-moved subject and no *by* phrase, passive with wh-moved object out of the *by* phrase, passive with wh-moved *by* phrase, passive with wh-moved indirect object and *by* phrase, passive with wh-moved indirect object and no *by* phrase, passive with relative clause on subject and *by* phrase with and without comp, passive with relative clause on subject and no *by* phrase with and without comp, passive with relative clause on object of the *by* phrase with and without comp/with PP

pied-piping, passive with relative clause on the indirect object and *by* phrase with and without comp/with PP pied-piping, passive with relative clause on the indirect object and no *by* phrase with and without comp/with PP pied-piping, passive with/without *by*-phrase with adjunct (gap-less) relative clause with comp/with PP pied-piping, gerund passive with *by* phrase, gerund passive without *by* phrase;

Shifted: wh-moved subject, wh-moved direct object, wh-moved object of PP, wh-moved PP, subject relative clause with and without comp, adjunct (gap-less) relative clause with comp/with PP pied-piping, direct object relative clause with comp/with PP pied-piping, object of PP relative clause with and without comp/with PP pied-piping, imperative, determiner gerund, NP gerund, passive with *by* phrase, passive without *by* phrase, passive with wh-moved subject and *by* phrase, passive with wh-moved subject and no *by* phrase, passive with wh-moved object out of the *by* phrase, passive with wh-moved *by* phrase, passive with wh-moved object out of the PP and *by* phrase, passive with wh-moved object out of the PP and no *by* phrase, passive with wh-moved PP and *by* phrase, passive with wh-moved PP and no *by* phrase, passive with relative clause on subject and *by* phrase with and without comp, passive with relative clause on subject and no *by* phrase with and without comp, passive with relative clause on object of the *by* phrase with and without comp/with PP pied-piping, passive with relative clause on the object of the PP and *by* phrase with and without comp/with PP pied-piping, passive with relative clause on the object of the PP and no *by* phrase with and without comp/with PP pied-piping, passive with/without *by*-phrase with adjunct (gap-less) relative clause with comp/with PP pied-piping, gerund passive with *by* phrase, gerund passive without *by* phrase.

6.6 Sentential Complement with NP: Tnx0Vnx1s2

Description: This tree family is selected by verbs that take both an NP and a sentential complement. The sentential complement may be infinitive or indicative. The type of clause is specified by each individual verb in its syntactic lexicon entry. A given verb may select more than one type of sentential complement. The declarative tree, and many other trees in this family, are auxiliary trees, as opposed to the more common initial trees. These auxiliary trees adjoin onto an S node in an existing tree of the type specified by the sentential complement. This is the mechanism by which TAGs are able to maintain long-distance dependencies (see Chapter 13), even over multiple embeddings (e.g. *What did Bill tell Mary that John said?*). 79 verbs select this tree family.

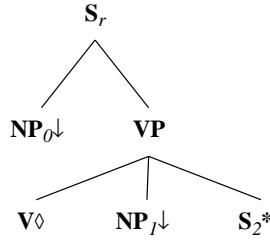
Examples: *beg, expect, tell*

Srini begged Mark to increase his disk quota .

Beth told Jim that it was his turn .

Declarative tree: See Figure 6.6.

Other available trees: wh-moved subject, wh-moved object, wh-moved sentential complement, subject relative clause with and without comp, adjunct (gap-less) relative clause with comp/with PP pied-piping, object relative clause with and without comp, imperative, determiner gerund, NP gerund, passive with *by* phrase before sentential complement, passive with *by* phrase after sentential complement, passive without *by* phrase, passive with

Figure 6.6: Declarative Sentential Complement with NP Tree: $\beta_{nx0}V_{nx1}s_2$

wh-moved subject and *by* phrase before sentential complement, passive with wh-moved subject and *by* phrase after sentential complement, passive with wh-moved subject and no *by* phrase, passive with wh-moved object out of the *by* phrase, passive with wh-moved *by* phrase, passive with relative clause on subject and *by* phrase before sentential complement with and without comp, passive with relative clause on subject and *by* phrase after sentential complement with and without comp, passive with relative clause on subject and no *by* phrase with and without comp, passive with/without *by*-phrase with adjunct (gap-less) relative clause with comp/with PP pied-piping, gerund passive with *by* phrase before sentential complement, gerund passive with *by* phrase after the sentential complement, gerund passive without *by* phrase, parenthetical reporting clause.

6.7 Intransitive Verb Particle: $T_{nx0}V_{pl}$

Description: The trees in this tree family are anchored by both the verb and the verb particle. Both appear in the syntactic lexicon and together select this tree family. Intransitive verb particles can be difficult to distinguish from intransitive verbs with adverbs adjoined on. The main diagnostics for including verbs in this class are whether the meaning is compositional or not, and whether there is a transitive version of the verb/verb particle combination with the same or similar meaning. The existence of an alternate compositional meaning is a strong indication for a separate verb particle construction. There are 159 verb/verb particle combinations.

Examples: *add up, come out, sign off*
The numbers never quite added up .
John finally came out (of the closet) .
I think that I will sign off now .

Declarative tree: See Figure 6.7.

Other available trees: wh-moved subject, subject relative clause with and without comp, adjunct (gap-less) relative clause with comp/with PP pied-piping, imperative, determiner gerund, NP gerund.

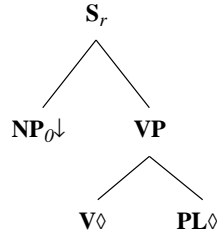


Figure 6.7: Declarative Intransitive Verb Particle Tree: $\alpha nx0Vpl$

6.8 Transitive Verb Particle: $Tnx0Vplnx1$

Description: Verb/verb particle combinations that take an NP complement select this tree family. Both the verb and the verb particle are anchors of the trees. Particle movement has been taken as the diagnostic to distinguish verb particle constructions from intransitives with adjoined PP's. If the alleged particle is able to undergo particle movement, in other words appear both before and after the direct object, then it is judged to be a particle. Items that do not undergo particle movement are taken to be prepositions. In many, but not all, of the verb particle cases, there is also an alternate prepositional meaning in which the lexical item did not move. (e.g. *He looked up the number (in the phonebook). He looked the number up. Srimi looked up the road (for Purnima's car). *He looked the road up.*) There are 489 verb/verb particle combinations.

Examples: *blow off, make up, pick out*

He blew off his linguistics class for the third time .

He blew his linguistics class off for the third time .

The dyslexic leprechaun made up the syntactic lexicon .

The dyslexic leprechaun made the syntactic lexicon up .

I would like to pick out a new computer .

I would like to pick a new computer out .

Declarative tree: See Figure 6.8.

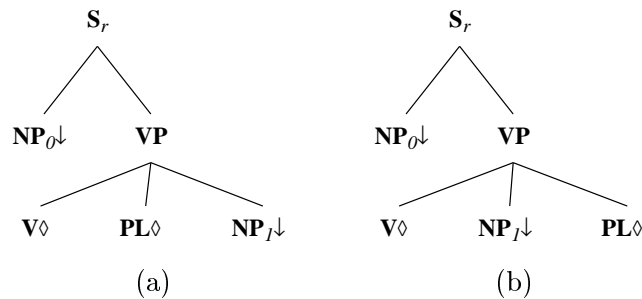


Figure 6.8: Declarative Transitive Verb Particle Tree: $\alpha nx0Vplnx1$ (a) and $\alpha nx0Vnx1pl$ (b)

Other available trees: wh-moved subject with particle before the NP, wh-moved subject with particle after the NP, wh-moved object, subject relative clause with particle before

the NP with and without comp, subject relative clause with particle after the NP with and without comp, object relative clause with and without comp, adjunct (gap-less) relative clause with particle before the NP with comp/with PP pied-piping, adjunct (gap-less) relative clause with particle after the NP with comp/with PP pied-piping, imperative with particle before the NP, imperative with particle after the NP, determiner gerund with particle before the NP, NP gerund with particle before the NP, NP gerund with particle after the NP, passive with *by* phrase, passive without *by* phrase, passive with wh-moved subject and *by* phrase, passive with wh-moved subject and no *by* phrase, passive with wh-moved object out of the *by* phrase, passive with wh-moved *by* phrase, passive with relative clause on subject and *by* phrase with and without comp, passive with relative clause on subject and no *by* phrase with and without comp, passive with relative clause on object of the *by* phrase with and without comp/with PP pied-piping, passive with/without *by*-phrase with adjunct (gap-less) relative clause with comp/with PP pied-piping, gerund passive with *by* phrase, gerund passive without *by* phrase.

6.9 Ditransitive Verb Particle: Tnx0Vplnx1nx2

Description: Verb/verb particle combinations that select this tree family take 2 NP complements. Both the verb and the verb particle anchor the trees, and the verb particle can occur before, between, or after the noun phrases. Perhaps because of the complexity of the sentence, these verbs do not seem to have passive alternations (**A new bank account was opened up Michelle by me*). There are 4 verb/verb particle combinations that select this tree family. The exhaustive list is given in the examples.

Examples: *dish out, open up, pay off, rustle up*
I opened up Michelle a new bank account .
I opened Michelle up a new bank account .
I opened Michelle a new bank account up .

Declarative tree: See Figure 6.9.

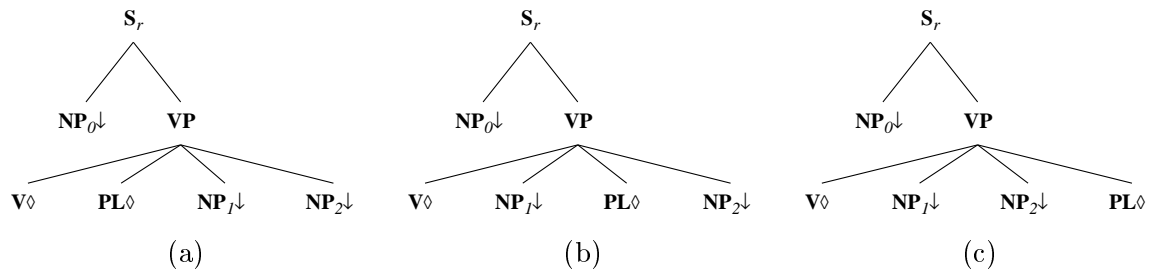


Figure 6.9: Declarative Ditransitive Verb Particle Tree: $\alpha nx0Vplnx1nx2$ (a), $\alpha nx0Vnx1plnx2$ (b) and $\alpha nx0Vnx1nx2pl$ (c)

Other available trees: wh-moved subject with particle before the NP's, wh-moved subject with particle between the NP's, wh-moved subject with particle after the NP's, wh-moved indirect object with particle before the NP's, wh-moved indirect object with particle after

the NP's, wh-moved direct object with particle before the NP's, wh-moved direct object with particle between the NP's, subject relative clause with particle before the NP's with and without comp, subject relative clause with particle between the NP's with and without comp, subject relative clause with particle after the NP's with and without comp, indirect object relative clause with particle before the NP's with and without comp, indirect object relative clause with particle after the NP's with and without comp, direct object relative clause with particle before the NP's with and without comp, direct object relative clause with particle between the NP's with and without comp, adjunct (gap-less) relative clause with comp/with PP pied-piping, imperative with particle before the NP's, imperative with particle between the NP's, imperative with particle after the NP's, determiner gerund with particle before the NP's, NP gerund with particle before the NP's, NP gerund with particle between the NP's, NP gerund with particle after the NP's.

6.10 Intransitive with PP: Tnx0Vpnx1

Description: The verbs that select this tree family are not strictly intransitive, in that they **must** be followed by a prepositional phrase. Verbs that are intransitive and simply **can** be followed by a prepositional phrase do not select this family, but instead have the PP adjoin onto the intransitive sentence. Accordingly, there should be no verbs in both this class and the intransitive tree family (see section 6.1). The prepositional phrase is not restricted to being headed by any particular lexical item. Note that these are not transitive verb particles (see section 6.8), since the head of the PP does not move. 169 verbs select this tree family.

Examples: *grab, impinge, provide*

Seth grabbed for the brass ring .

The noise gradually impinged on Dania's thoughts .

A good host provides for everyone's needs .

Declarative tree: See Figure 6.10.

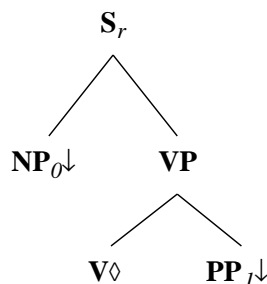


Figure 6.10: Declarative Intransitive with PP Tree: $\alpha\text{nx0Vpnx1}$

Other available trees: wh-moved subject, wh-moved object of the PP, wh-moved PP, subject relative clause with and without comp, adjunct (gap-less) relative clause with comp/with

PP pied-piping, object of the PP relative clause with and without comp/with PP pied-piping, imperative, determiner gerund, NP gerund, passive with *by* phrase, passive without *by* phrase, passive with wh-moved subject and *by* phrase, passive with wh-moved subject and no *by* phrase, passive with wh-moved *by* phrase, passive with relative clause on subject and *by* phrase with and without comp, passive with relative clause on subject and no *by* phrase with and without comp, passive with relative clause on object of the *by* phrase with and without comp/with PP pied-piping, passive with/without *by*-phrase with adjunct (gap-less) relative clause with comp/with PP pied-piping, gerund passive with *by* phrase, gerund passive without *by* phrase.

6.11 Predicative Multi-word with Verb, Prep anchors: Tnx0VPnx1

Description: This tree family is selected by multiple anchor verb/preposition pairs which together have a non-compositional interpretation. For example, *think of* has the non-compositional interpretation involving the inception of a notion or mental entity in addition to the interpretation in which the agent is thinking about someone or something. Anchors for this tree must be able to take both gerunds and regular NP's in the second noun position. To allow adverbs to appear between the verb and the preposition, the trees contain an extra VP level. Several of the verbs which select the Tnx0Vpnx1 family, but which should not have quite the freedom it allows, will be moving to this family for the next release. 28 verb/preposition pairs select this tree family.

Examples: *think of, believe in, depend on*
Calvin thought of a new idea .
Hobbes believes in sleeping all day .
Bill depends on drinking coffee for stimulation .

Declarative tree: See Figure 6.11.

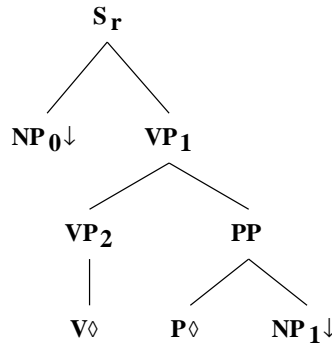


Figure 6.11: Declarative PP Complement Tree: α nx0VPnx1

Other available trees: wh-moved subject, wh-moved object, subject relative clause with and without comp, adjunct (gap-less) relative clause with comp/with PP pied-piping, object relative clause with and without comp, imperative, determiner gerund, NP gerund, passive

with *by* phrase, passive without *by* phrase, passive with wh-moved subject and *by* phrase, passive with wh-moved subject and no *by* phrase, passive with wh-moved object out of the *by* phrase, passive with wh-moved *by* phrase, passive with relative clause on subject and *by* phrase with and without comp, passive with relative clause on subject and no *by* phrase with and without comp, passive with relative clause on object on the *by* phrase with and without comp/with PP pied-piping, passive with/without *by*-phrase with adjunct (gap-less) relative clause with comp/with PP pied-piping, gerund passive with *by* phrase, gerund passive without *by* phrase. In addition, two other trees that allow transitive verbs to function as adjectives (e.g. *the thought of idea*) are also in the family.

6.12 Sentential Complement: Tnx0Vs1

Description: This tree family is selected by verbs that take just a sentential complement. The sentential complement may be of type infinitive, indicative, or small clause (see Chapter 9). The type of clause is specified by each individual verb in its syntactic lexicon entry, and a given verb may select more than one type of sentential complement. The declarative tree, and many other trees in this family, are auxiliary trees, as opposed to the more common initial trees. These auxiliary trees adjoin onto an S node in an existing tree of the type specified by the sentential complement. This is the mechanism by which TAGs are able to maintain long-distance dependencies (see Chapter 13), even over multiple embeddings (e.g. *What did Bill think that John said?*). 338 verbs select this tree family.

Examples: *consider, think*

Dania considered the algorithm unworkable .

Srini thought that the program was working .

Declarative tree: See Figure 6.12.

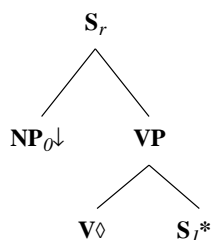


Figure 6.12: Declarative Sentential Complement Tree: β nx0Vs1

Other available trees: wh-moved subject, wh-moved sentential complement, subject relative clause with and without comp, adjunct (gap-less) relative clause with comp/with PP pied-piping, imperative, determiner gerund, NP gerund, parenthetical reporting clause.

6.13 Intransitive with Adjective: Tnx0Vax1

Description: The verbs that select this tree family take an adjective as a complement. The

adjective may be regular, comparative, or superlative. It may also be formed from the special class of adjectives derived from the transitive verbs (e.g. *agitated, broken*). See section 6.2). Unlike the Intransitive with PP verbs (see section 6.10), some of these verbs may also occur as bare intransitives as well. This distinction is drawn because adjectives do not normally adjoin onto sentences, as prepositional phrases do. Other intransitive verbs can only occur with the adjective, and these select only this family. The verb class is also distinguished from the adjective small clauses (see section 6.20) because these verbs are not raising verbs. 34 verbs select this tree family.

Examples: *become, grow, smell*

The greenhouse became hotter .

The plants grew tall and strong .

The flowers smelled wonderful .

Declarative tree: See Figure 6.13.

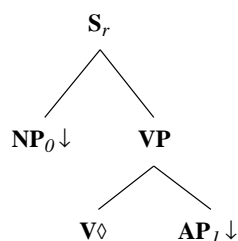


Figure 6.13: Declarative Intransitive with Adjective Tree: $\alpha nx0Vax1$

Other available trees: wh-moved subject, wh-moved adjective (*how*), subject relative clause with and without comp, adjunct (gap-less) relative clause with comp/with PP pied-piping, imperative, NP gerund.

6.14 Transitive Sentential Subject: Ts0Vnx1

Description: The verbs that select this tree family all take sentential subjects, and are often referred to as ‘psych’ verbs, since they all refer to some psychological state of mind. The sentential subject can be indicative (complementizer required) or infinitive (complementizer optional). 100 verbs that select this tree family.

Examples: *delight, impress, surprise*

that the tea had rosehips in it delighted Christy .

to even attempt a marathon impressed Dania .

For Jim to have walked the dogs surprised Beth .

Declarative tree: See Figure 6.14.

Other available trees: wh-moved subject, wh-moved object, subject relative clause with and without comp, adjunct (gap-less) relative clause with comp/with PP pied-piping.

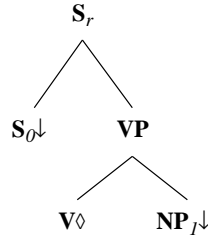


Figure 6.14: Declarative Sentential Subject Tree: $\alpha s0Vnx1$

6.15 Light Verbs: Tnx0lVN1

Description: The verb/noun pairs that select this tree families are pairs in which the interpretation is non-compositional and the noun contributes argument structure to the predicate (e.g. *The man took a walk.* vs. *The man took a radio*). The verb and the noun occur together in the syntactic database, and both anchor the trees. The verbs in the light verb constructions are *do*, *give*, *have*, *make* and *take*. The noun following the light verb is (usually) in a bare infinitive form (*have a good cry*) and usually occurs with *a(n)*. However, we include deverbal nominals (*take a bath*, *give a demonstration*) as well. Constructions with nouns that do not contribute an argument structure (*have a cigarette*, *give NP a black eye*) are excluded. In addition to semantic considerations of light verbs, they differ syntactically from Transitive verbs (section 6.2) as well in that the noun in the light verb construction does not extract. Some of the verb-noun anchors for this family, like *take aim* and *take hold* disallow determiners, while others require particular determiners. For example, *have think* must be indefinite and singular, as attested by the ungrammaticality of **John had the think/some thinks*. Another anchor, *take leave* can occur either bare or with a possessive pronoun (e.g., *John took his leave*, but not **John took the leave*). This is accomplished through feature specification on the lexical entries. There are 259 verb/noun pairs that select the light verb tree.

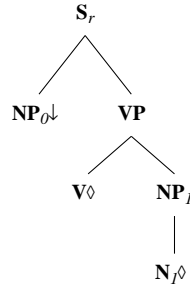
Examples: *give groan*, *have discussion*, *make comment*
The audience gave a collective groan .
We had a big discussion about closing the libraries .
The professors made comments on the paper .

Declarative tree: See Figure 6.15.

Other available trees: wh-moved subject, subject relative clause with and without comp, adjunct (gap-less) relative clause with comp/with PP pied-piping, imperative, determiner gerund, NP gerund.

6.16 Ditransitive Light Verbs with PP Shift: Tnx0lVN1Pnx2

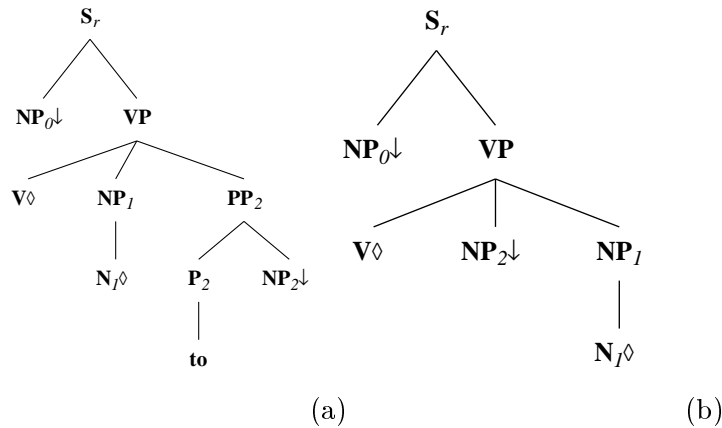
Description: The verb/noun pairs that select this tree family are pairs in which the interpretation is non-compositional and the noun contributes argument structure to the predicate (e.g. *Dania made Srini a cake.* vs. *Dania made Srini a loan*). The verb and the noun


 Figure 6.15: Declarative Light Verb Tree: $\alpha nx0lVN1$

occur together in the syntactic database, and both anchor the trees. The verbs in these light verb constructions are *give* and *make*. The noun following the light verb is (usually) a bare infinitive form (e.g. *make a promise to Anoop*). However, we include deverbal nominals (e.g. *make a payment to Anoop*) as well. Constructions with nouns that do not contribute an argument structure are excluded. In addition to semantic considerations of light verbs, they differ syntactically from the Ditransitive with PP Shift verbs (see section 6.5) as well in that the noun in the light verb construction does not extract. Also, passivization is severely restricted. Special determiner requirements and restrictions are handled in the same manner as for the $Tnx0lVN1$ family. There are 18 verb/noun pairs that select this family.

Examples: *give look, give wave, make promise*
Dania gave Carl a murderous look .
Amanda gave us a little wave as she left .
Dania made Doug a promise .

Declarative tree: See Figure 6.16.


 Figure 6.16: Declarative Light Verbs with PP Tree: $\alpha nx0lVN1Pnx2$ (a), $\alpha nx0lVnx2N1$ (b)

Other available trees: Non-shifted: wh-moved subject, wh-moved indirect object, subject relative clause with and without comp, adjunct (gap-less) relative clause with comp/with PP pied-piping, indirect object relative clause with and without comp/with PP pied-piping, imperative, NP gerund, passive with *by* phrase, passive with *by*-phrase with adjunct (gap-less) relative clause with comp/with PP pied-piping, gerund passive with *by* phrase, gerund passive without *by* phrase

Shifted: wh-moved subject, wh-moved object of PP, wh-moved PP, subject relative clause with and without comp, object of PP relative clause with and without comp/with PP pied-piping, imperative, determiner gerund, NP gerund, passive with *by* phrase with adjunct (gap-less) relative clause with comp/with PP pied-piping, gerund passive with *by* phrase, gerund passive without *by* phrase.

6.17 NP It-Cleft: TItVnx1s2

Description: This tree family is selected by *be* as the main verb and *it* as the subject. Together these two items serve as a multi-component anchor for the tree family. This tree family is used for it-clefts in which the clefted element is an NP and there are no gaps in the clause which follows the NP. The NP is interpreted as an adjunct of the following clause. See Chapter 11 for additional discussion.

Examples: *it be*

it was yesterday that we had the meeting .

Declarative tree: See Figure 6.17.

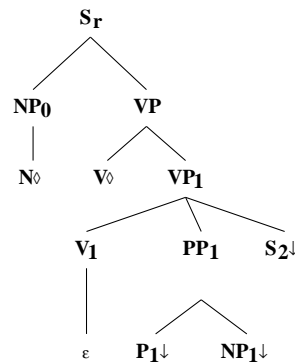


Figure 6.17: Declarative NP It-Cleft Tree: α ItVpnx1s2

Other available trees: inverted question, wh-moved object with *be* inverted, wh-moved object with *be* not inverted, adjunct (gap-less) relative clause with comp/with PP pied-piping.

6.18 PP It-Cleft: TItVpnx1s2

Description: This tree family is selected by *be* as the main verb and *it* as the subject. Together these two items serve as a multi-component anchor for the tree family. This tree family is

used for it-clefts in which the clefted element is a PP and there are no gaps in the clause which follows the PP. The PP is interpreted as an adjunct of the following clause. See Chapter 11 for additional discussion.

Examples: *it be*

it was at Kent State that the police shot all those students .

Declarative tree: See Figure 6.18.

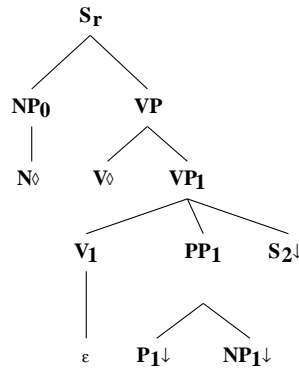


Figure 6.18: Declarative PP It-Cleft Tree: $\alpha\text{ItVnx1s2}$

Other available trees: inverted question, wh-moved prepositional phrase with *be* inverted, wh-moved prepositional phrase with *be* not inverted, adjunct (gap-less) relative clause with comp/with PP pied-piping.

6.19 Adverb It-Cleft: TItVad1s2

Description: This tree family is selected by *be* as the main verb and *it* as the subject. Together these two items serve as a multi-component anchor for the tree family. This tree family is used for it-clefts in which the clefted element is an adverb and there are no gaps in the clause which follows the adverb. The adverb is interpreted as an adjunct of the following clause. See Chapter 11 for additional discussion.

Examples: *it be*

it was reluctantly that Dania agreed to do the tech report .

Declarative tree: See Figure 6.19.

Other available trees: inverted question, wh-moved adverb *how* with *be* inverted, wh-moved adverb *how* with *be* not inverted, adjunct (gap-less) relative clause with comp/with PP pied-piping.

6.20 Adjective Small Clause Tree: Tnx0Ax1

Description: These trees are not anchored by verbs, but by adjectives. They are explained

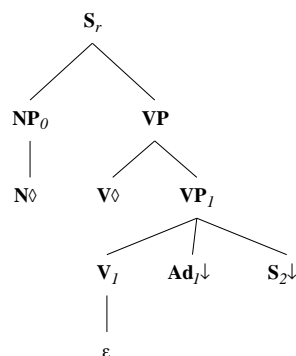


Figure 6.19: Declarative Adverb It-Cleft Tree: $\alpha\text{ItVad1s2}$

in much greater detail in the section on small clauses (see section 9.3). This section is presented here for completeness. 3244 adjectives select this tree family.

Examples: *addictive, dangerous, wary*
cigarettes are addictive .
smoking cigarettes is dangerous .
John seems wary of the Surgeon General's warnings .

Declarative tree: See Figure 6.20.

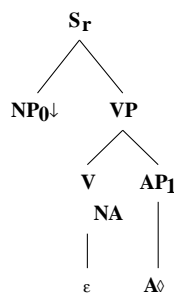


Figure 6.20: Declarative Adjective Small Clause Tree: αnx0Ax1

Other available trees: wh-moved subject, wh-moved adjective *how*, relative clause on subject with and without comp, imperative, NP gerund, adjunct (gap-less) relative clause with comp/with PP pied-piping.

6.21 Adjective Small Clause with Sentential Complement: Tnx0A1s1

Description: This tree family is selected by adjectives that take sentential complements. The sentential complements can be indicative or infinitive. Note that these trees are anchored by adjectives, not verbs. Small clauses are explained in much greater detail in section 9.3. This section is presented here for completeness. 669 adjectives select this tree family.

Examples: *able, curious, disappointed*

Christy was able to find the problem .

Christy was curious whether the new analysis was working .

Christy was sad that the old analysis failed .

Declarative tree: See Figure 6.21.

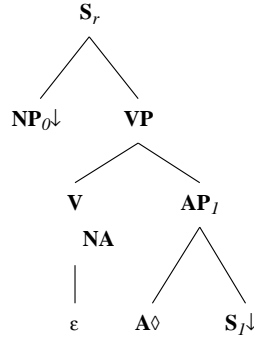


Figure 6.21: Declarative Adjective Small Clause with Sentential Complement Tree: $\alpha n x 0 A 1 s 1$

Other available trees: wh-moved subject, wh-moved adjective *how*, relative clause on subject with and without comp, imperative, NP gerund, adjunct (gap-less) relative clause with comp/with PP pied-piping.

6.22 Adjective Small Clause with Sentential Subject: $Ts 0 A x 1$

Description: This tree family is selected by adjectives that take sentential subjects. The sentential subjects can be indicative or infinitive. Note that these trees are anchored by adjectives, not verbs. Most adjectives that take the Adjective Small Clause tree family (see section 6.20) take this family as well.⁶ Small clauses are explained in much greater detail in section 9.3. This section is presented here for completeness. 3,185 adjectives select this tree family.

Examples: *decadent, incredible, uncertain*

to eat raspberry chocolate truffle ice cream is decadent .

that Carl could eat a large bowl of it is incredible .

whether he will actually survive the experience is uncertain .

Declarative tree: See Figure 6.22.

Other available trees: wh-moved subject, wh-moved adjective, adjunct (gap-less) relative clause with comp/with PP pied-piping.

⁶No great attempt has been made to go through and decide which adjectives should actually take this family and which should not.

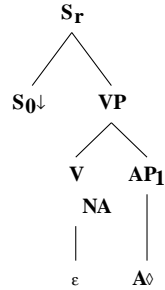


Figure 6.22: Declarative Adjective Small Clause with Sentential Subject Tree: $\alpha s_0 A x_1$

6.23 Equative *BE*: $T_{nx_0 B E n x_1}$

Description: This tree family is selected only by the verb *be*. It is distinguished from the predicative NP's (see section 6.24) in that two NP's are equated, and hence interchangeable (see Chapter 9 for more discussion on the English copula and predicative sentences). The XTAG analysis for equative *be* is explained in greater detail in section 9.4.

Examples: *be*

That man is my uncle.

Declarative tree: See Figure 6.23.

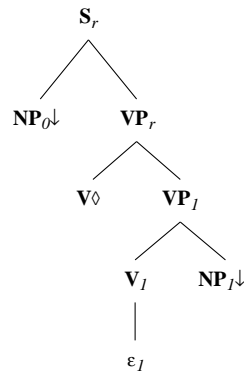


Figure 6.23: Declarative Equative *BE* Tree: $\alpha n x_0 B E n x_1$

Other available trees: inverted-question.

6.24 NP Small Clause: $T_{nx_0 N 1}$

Description: The trees in this tree family are not anchored by verbs, but by nouns. Small clauses are explained in much greater detail in section 9.3. This section is presented here for completeness. 5,595 nouns select this tree family.

Examples: *author, chair, dish*

Dania is an author .

that blue, warped-looking thing is a chair .

those broken pieces were dishes .

Declarative tree: See Figure 6.24.

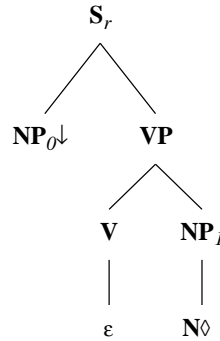


Figure 6.24: Declarative NP Small Clause Trees: $\alpha n x 0 N 1$

Other available trees: wh-moved subject, wh-moved object, relative clause on subject with and without comp, imperative, NP gerund, adjunct (gap-less) relative clause with comp/with PP pied-piping.

6.25 NP Small Clause with Sentential Complement: $T n x 0 N 1 s 1$

Description: This tree family is selected by the small group of nouns that take sentential complements by themselves (see section 8.8). The sentential complements can be indicative or infinitive, depending on the noun. Small clauses in general are explained in much greater detail in the section 9.3. This section is presented here for completeness. 141 nouns select this family.

Examples: *admission, claim, vow*

The affidavits are admissions that they killed the sheep .

there is always the claim that they were insane .

this is his vow to fight the charges .

Declarative tree: See Figure 6.25.

Other available trees: wh-moved subject, wh-moved object, relative clause on subject with and without comp, imperative, NP gerund, adjunct (gap-less) relative clause with comp/with PP pied-piping.

6.26 NP Small Clause with Sentential Subject: $T s 0 N 1$

Description: This tree family is selected by nouns that take sentential subjects. The sentential subjects can be indicative or infinitive. Note that these trees are anchored by nouns, not

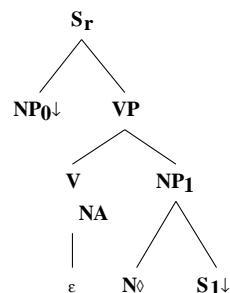


Figure 6.25: Declarative NP with Sentential Complement Small Clause Tree: $\alpha n x 0 N 1 s 1$

verbs. Most nouns that take the NP Small Clause tree family (see section 6.24) take this family as well.⁷ Small clauses are explained in much greater detail in section 9.3. This section is presented here for completeness. 5,519 nouns select this tree family.

Examples: *dilemma, insanity, tragedy*
whether to keep the job he hates is a dilemma .
to invest all of your money in worms is insanity .
that the worms died is a tragedy .

Declarative tree: See Figure 6.26.

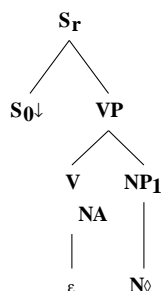


Figure 6.26: Declarative NP Small Clause with Sentential Subject Tree: $\alpha s 0 N 1$

Other available trees: wh-moved subject, adjunct (gap-less) relative clause with comp/with PP pied-piping.

6.27 PP Small Clause: $T n x 0 P n x 1$

Description: This family is selected by prepositions that can occur in small clause constructions. For more information on small clause constructions, see section 9.3. This section is presented here for completeness. 39 prepositions select this tree family.

⁷No great attempt has been made to go through and decide which nouns should actually take this family and which should not.

Examples: *around, in, underneath*

Chris is around the corner .

Trisha is in big trouble .

The dog is underneath the table .

Declarative tree: See Figure 6.27.

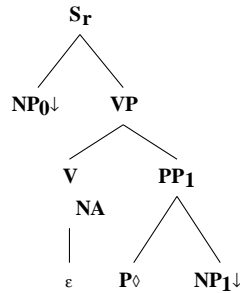


Figure 6.27: Declarative PP Small Clause Tree: $\alpha n x 0 P n x 1$

Other available trees: wh-moved subject, wh-moved object of PP, relative clause on subject with and without comp, relative clause on object of PP with and without comp/with PP pied-piping, imperative, NP gerund, adjunct (gap-less) relative clause with comp/with PP pied-piping.

6.28 Exhaustive PP Small Clause: $T n x 0 P x 1$

Description: This family is selected by **exhaustive** prepositions that can occur in small clauses. Exhaustive prepositions are prepositions that function as prepositional phrases by themselves. For more information on small clause constructions, please see section 9.3. The section is included here for completeness. 33 exhaustive prepositions select this tree family.

Examples: *abroad, below, outside*

Dr. Joshi is abroad .

The workers are all below .

Clove is outside .

Declarative tree: See Figure 6.28.

Other available trees: wh-moved subject, wh-moved PP, relative clause on subject with and without comp, imperative, NP gerund, adjunct (gap-less) relative clause with comp/with PP pied-piping.

6.29 PP Small Clause with Sentential Subject: $T s 0 P n x 1$

Description: This tree family is selected by prepositions that take sentential subjects. The sentential subject can be indicative or infinitive. Small clauses are explained in much

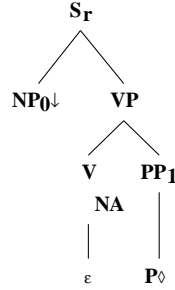


Figure 6.28: Declarative Exhaustive PP Small Clause Tree: $\alpha nx0Px1$

greater detail in section 9.3. This section is presented here for completeness. 39 prepositions select this tree family.

Examples: *beyond, unlike*

that Ken could forget to pay the taxes is beyond belief .

to explain how this happened is outside the scope of this discussion .

for Ken to do something right is unlike him .

Declarative tree: See Figure 6.29.

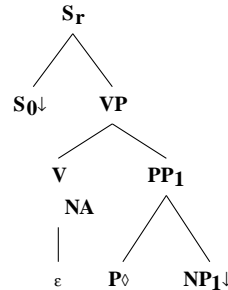


Figure 6.29: Declarative PP Small Clause with Sentential Subject Tree: $\alpha s0Pnx1$

Other available trees: wh-moved subject, relative clause on object of the PP with and without comp/with PP pied-piping, adjunct (gap-less) relative clause with comp/with PP pied-piping.

6.30 Intransitive Sentential Subject: Ts0V

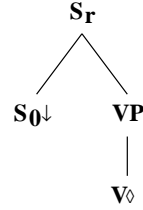
Description: Only the verb *matter* selects this tree family. The sentential subject can be indicative (complementizer required) or infinitive (complementizer optional).

Examples: *matter*

to arrive on time matters considerably .

that Joshi attends the meetings matters to everyone .

Declarative tree: See Figure 6.30.

Figure 6.30: Declarative Intransitive Sentential Subject Tree: $\alpha s0V$

Other available trees: wh-moved subject, adjunct (gap-less) relative clause with comp/with PP pied-piping.

6.31 Sentential Subject with ‘to’ complement: $Ts0Vtonx1$

Description: The verbs that select this tree family are *fall*, *occur* and *leak*. The sentential subject can be indicative (complementizer required) or infinitive (complementizer optional).

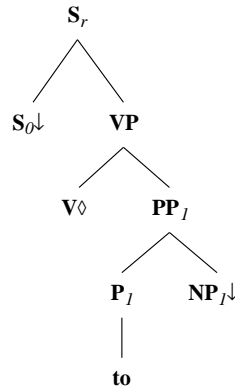
Examples: *fall*, *occur*, *leak*

to wash the car fell to the children .

that he should leave occurred to the party crasher .

whether the princess divorced the prince leaked to the press .

Declarative tree: See Figure 6.31.

Figure 6.31: Sentential Subject Tree with ‘to’ complement: $\alpha s0Vtonx1$

Other available trees: wh-moved subject, adjunct (gap-less) relative clause with comp/with PP pied-piping.

6.32 PP Small Clause, with Adv and Prep anchors: $Tnx0ARBPnx1$

Description: This family is selected by multi-word prepositions that can occur in small clause constructions. In particular, this family is selected by two-word prepositions, where the

first word is an adverb, the second word a preposition. Both components of the multi-word preposition are anchors. For more information on small clause constructions, see section 9.3. 8 multi-word prepositions select this tree family.

Examples: *ahead of, close to*

The little girl is ahead of everyone else in the race .

The project is close to completion .

Declarative tree: See Figure 6.32.

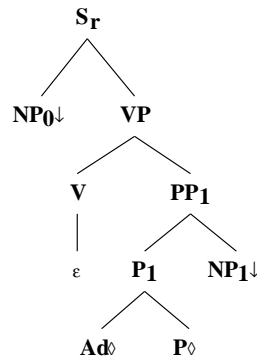


Figure 6.32: Declarative PP Small Clause tree with two-word preposition, where the first word is an adverb, and the second word is a preposition: $\alpha x0ARBPnx1$

Other available trees: wh-moved subject, wh-moved object of PP, relative clause on subject with and without comp, relative clause on object of PP with and without comp, adjunct (gap-less) relative clause with comp/with PP pied-piping, imperative, NP Gerund.

6.33 PP Small Clause, with Adj and Prep anchors: $Tnx0APnx1$

Description: This family is selected by multi-word prepositions that can occur in small clause constructions. In particular, this family is selected by two-word prepositions, where the first word is an adjective, the second word a preposition. Both components of the multi-word preposition are anchors. For more information on small clause constructions, see section 9.3. 8 multi-word prepositions select this tree family.

Examples: *according to, void of*

The operation we performed was according to standard procedure .

He is void of all feeling .

Declarative tree: See Figure 6.33.

Other available trees: wh-moved subject, relative clause on subject with and without comp, relative clause on object of PP with and without comp, wh-moved object of PP, adjunct (gap-less) relative clause with comp/with PP pied-piping.

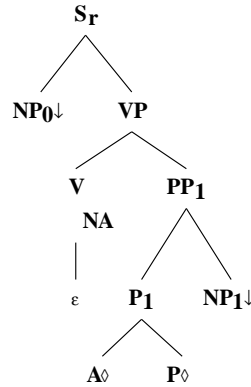


Figure 6.33: Declarative PP Small Clause tree with two-word preposition, where the first word is an adjective, and the second word is a preposition: $\alpha nx0APnx1$

6.34 PP Small Clause, with Noun and Prep anchors: $Tnx0NPnx1$

Description: This family is selected by multi-word prepositions that can occur in small clause constructions. In particular, this family is selected by two-word prepositions, where the first word is a noun, the second word a preposition. Both components of the multi-word preposition are anchors. For more information on small clause constructions, see section 9.3. 1 multi-word preposition selects this tree family.

Examples: *thanks to*

The fact that we are here tonight is thanks to the valiant efforts of our staff .

Declarative tree: See Figure 6.34.

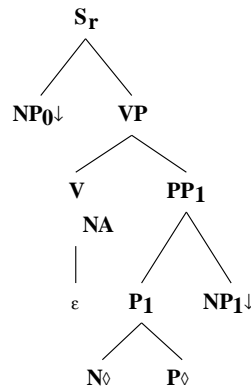


Figure 6.34: Declarative PP Small Clause tree with two-word preposition, where the first word is a noun, and the second word is a preposition: $\alpha nx0NPnx1$

Other available trees: wh-moved subject, wh-moved object of PP, relative clause on subject with and without comp, relative clause on object with comp, adjunct (gap-less) relative

clause with comp/with PP pied-piping.

6.35 PP Small Clause, with Prep anchors: Tnx0PPnx1

Description: This family is selected by multi-word prepositions that can occur in small clause constructions. In particular, this family is selected by two-word prepositions, where both words are prepositions. Both components of the multi-word preposition are anchors. For more information on small clause constructions, see section 9.3. 9 multi-word prepositions select this tree family.

Examples: *on to, inside of*
that detective is on to you .
The red box is inside of the blue box .

Declarative tree: See Figure 6.35.

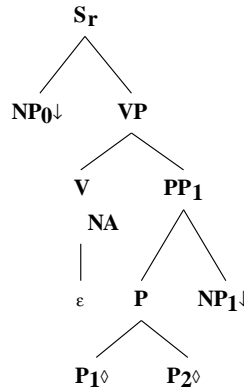


Figure 6.35: Declarative PP Small Clause tree with two-word preposition, where both words are prepositions: $\alpha nx0PPnx1$

Other available trees: wh-moved subject, wh-moved object of PP, relative clause on subject with and without comp, relative clause on object of PP with and without comp/with PP pied-piping, imperative, wh-moved object of PP, adjunct (gap-less) relative clause with comp/with PP pied-piping.

6.36 PP Small Clause, with Prep and Noun anchors: Tnx0PNaPnx1

Description: This family is selected by multi-word prepositions that can occur in small clause constructions. In particular, this family is selected by three-word prepositions. The first and third words are always prepositions, and the middle word is a noun. The noun is marked for null adjunction since it cannot be modified by noun modifiers. All three components of the multi-word preposition are anchors. For more information on small clause constructions, see section 9.3. 24 multi-word preposition select this tree family.

Examples: *in back of, in line with, on top of*

The red plaid box should be in back of the plain black box .

The evidence is in line with my newly concocted theory .

She is on top of the world .

**She is on direct top of the world .*

Declarative tree: See Figure 6.36.

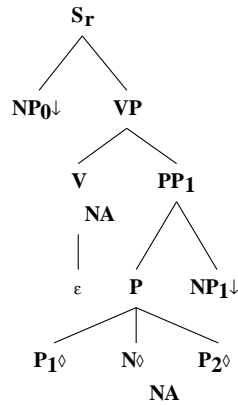


Figure 6.36: Declarative PP Small Clause tree with three-word preposition, where the middle noun is marked for null adjunction: $\alpha n x 0 P N a P n x 1$

Other available trees: wh-moved subject, wh-moved object of PP, relative clause on subject with and without comp, relative clause on object of PP with and without comp/with PP pied-piping, adjunct (gap-less) relative clause with comp/with PP pied-piping, imperative, NP Gerund.

6.37 PP Small Clause with Sentential Subject, and Adv and Prep anchors: Ts0ARBPNx1

Description: This tree family is selected by multi-word prepositions that take sentential subjects. In particular, this family is selected by two-word prepositions, where the first word is an adverb, the second word a preposition. Both components of the multi-word preposition are anchors. The sentential subject can be indicative or infinitive. Small clauses are explained in much greater detail in section 9.3. 2 prepositions select this tree family.

Examples: *due to, contrary to*

that David slept until noon is due to the fact that he never sleeps during the week .

that Michael's joke was funny is contrary to the usual status of his comic attempts .

Declarative tree: See Figure 6.37.

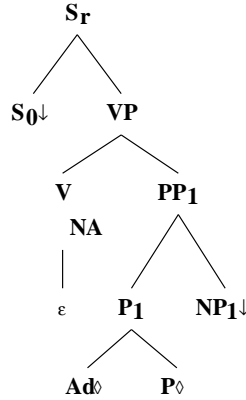


Figure 6.37: Declarative PP Small Clause with Sentential Subject Tree, with two-word preposition, where the first word is an adverb, and the second word is a preposition: $\alpha s0ARBPNx1$

Other available trees: wh-moved subject, relative clause on object of the PP with and without comp, adjunct (gap-less) relative clause with comp/with PP pied-piping.

6.38 PP Small Clause with Sentential Subject, and Adj and Prep anchors: $Ts0APnx1$

Description: This tree family is selected by multi-word prepositions that take sentential subjects. In particular, this family is selected by two-word prepositions, where the first word is an adjective, the second word a preposition. Both components of the multi-word preposition are anchors. The sentential subject can be indicative or infinitive. Small clauses are explained in much greater detail in section 9.3. 5 prepositions select this tree family.

Examples: *devoid of, according to*
that he could walk out on her is devoid of all reason .
that the conversation erupted precisely at that moment was according to my theory .

Declarative tree: See Figure 6.38.

Other available trees: wh-moved subject, relative clause on object of the PP with and without comp, adjunct (gap-less) relative clause with comp/with PP pied-piping.

6.39 PP Small Clause with Sentential Subject, and Noun and Prep anchors: $Ts0NPNx1$

Description: This tree family is selected by multi-word prepositions that take sentential subjects. In particular, this family is selected by two-word prepositions, where the first word is a noun, the second word a preposition. Both components of the multi-word preposi-

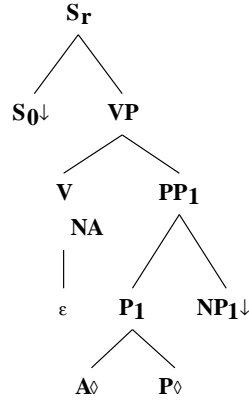


Figure 6.38: Declarative PP Small Clause with Sentential Subject Tree, with two-word preposition, where the first word is an adjective, and the second word is a preposition: $\alpha s0AP_{nx1}$

tion are anchors. The sentential subject can be indicative or infinitive. Small clauses are explained in much greater detail in section 9.3. 1 preposition selects this tree family.

Examples: *thanks to*

that she is worn out is thanks to a long day in front of the computer terminal .

Declarative tree: See Figure 6.39.

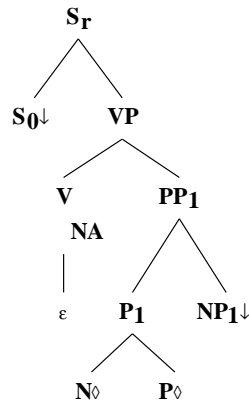


Figure 6.39: Declarative PP Small Clause with Sentential Subject Tree, with two-word preposition, where the first word is a noun, and the second word is a preposition: $\alpha s0NP_{nx1}$

Other available trees: wh-moved subject, relative clause on object of the PP with and without comp, adjunct (gap-less) relative clause with comp/with PP pied-piping.

6.40 PP Small Clause with Sentential Subject, and Prep anchors: Ts0PPnx1

Description: This tree family is selected by multi-word prepositions that take sentential subjects. In particular, this family is selected by two-word prepositions, where both words are prepositions. Both components of the multi-word preposition are anchors. The sentential subject can be indicative or infinitive. Small clauses are explained in much greater detail in section 9.3. 3 prepositions select this tree family.

Examples: *outside of*

that Mary did not complete the task on time is outside of the scope of this discussion .

Declarative tree: See Figure 6.40.

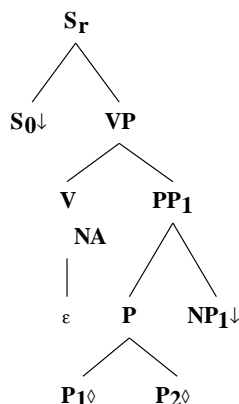


Figure 6.40: Declarative PP Small Clause with Sentential Subject Tree, with two-word preposition, where both words are prepositions: αs0PPnx1

Other available trees: wh-moved subject, relative clause on object of the PP with and without comp, adjunct (gap-less) relative clause with comp/with PP pied-piping.

6.41 PP Small Clause with Sentential Subject, and Prep and Noun anchors: Ts0PNaPnx1

Description: This tree family is selected by multi-word prepositions that take sentential subjects. In particular, this family is selected by three-word prepositions. The first and third words are always prepositions, and the middle word is a noun. The noun is marked for null adjunction since it cannot be modified by noun modifiers. All three components of the multi-word preposition are anchors. Small clauses are explained in much greater detail in section 9.3. 9 prepositions select this tree family.

Examples: *on account of, in support of*

that Joe had to leave the beach was on account of the hurricane .

that Maria could not come is in support of my theory about her .

**that Maria could not come is in direct/strict/desparate support of my theory about her .*

Declarative tree: See Figure 6.41.

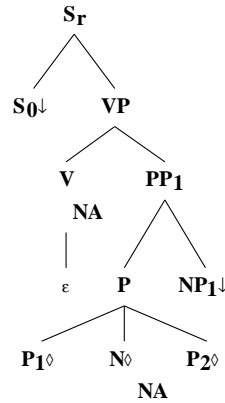


Figure 6.41: Declarative PP Small Clause with Sentential Subject Tree, with three-word preposition, where the middle noun is marked for null adjunction: $\alpha s0PNaPnx1$

Other available trees: wh-moved subject, relative clause on object of the PP with and without comp, adjunct (gap-less) relative clause with comp/with PP pied-piping.

6.42 Predicative Adjective with Sentential Subject and Complement: Ts0A1s1

Description: This tree family is selected by predicative adjectives that take sentential subjects and a sentential complement. This tree family is selected by *likely* and *certain*.

Examples: *likely, certain*

that Max continues to drive a Jaguar is certain to make Bill jealous .

for the Jaguar to be towed seems likely to make Max very angry .

Declarative tree: See Figure 6.42.

Other available trees: wh-moved subject, adjunct (gap-less) relative clause with comp/with PP pied-piping.

6.43 Locative Small Clause with Ad anchor: Tnx0nx1ARB

Description: These trees are not anchored by verbs, but by adverbs that are part of locative adverbial phrases. Locatives are explained in much greater detail in the section on the locative modifier trees (see section 19.6). The only remarkable aspect of this tree family is

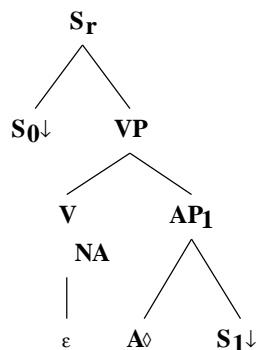


Figure 6.42: Predicative Adjective with Sentential Subject and Complement: $\alpha s0A1s1$

the wh-moved locative tree, $\alpha W1nx0nx1ARB$, shown in Figure 6.44. This is the only tree family with this type of transformation, in which the entire adverbial phrase is wh-moved but not all elements are replaced by wh items (as in *how many city blocks away is the record store?*). Locatives that consist of just the locative adverb or the locative adverb and a degree adverb (see Section 19.6 for details) are treated as exhaustive PPs and therefore select that tree family (Section 6.28) when used predicatively. For an extensive description of small clauses, see Section 9.3. 26 adverbs select this tree family.

Examples: *ahead, offshore, behind*
the crash is three blocks ahead
the naval battle was many kilometers offshore
how many blocks behind was Max?

Declarative tree: See Figure 6.43.

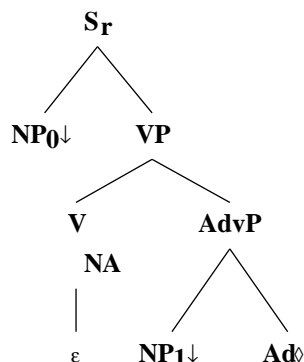
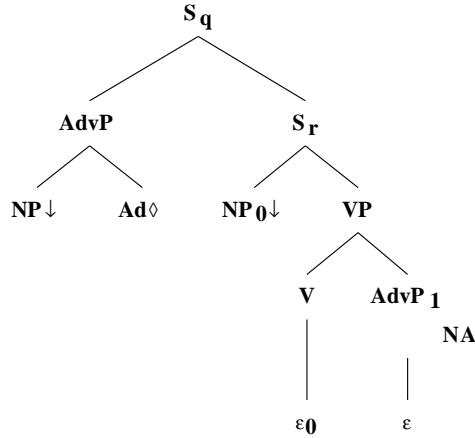


Figure 6.43: Declarative Locative Adverbial Small Clause Tree: $\alpha nx0nx1ARB$

Other available trees: wh-moved subject, relative clause on subject with and without comp, wh-moved locative, imperative, NP gerund.

Figure 6.44: Wh-moved Locative Small Clause Tree: $\alpha W1nx0nx1ARB$

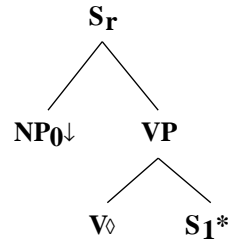
6.44 Exceptional Case Marking: $TXnx0Vs1$

Description: This tree family is selected by verbs that are classified as exceptional case marking, meaning that the verb assigns accusative case to the subject of the sentential complement. This is in contrast to verbs in the $Tnx0Vnx1s2$ family (section 6.6), which assign accusative case to a NP which is not part of the sentential complement. ECM verbs take sentential complements which are either an infinitive or a “bare” infinitive. As with the $Tnx0Vs1$ family (section 6.12), the declarative and other trees in the $Xnx0Vs1$ family are auxiliary trees, as opposed to the more common initial trees. These auxiliary trees adjoin onto an S node in an existing tree of the type specified by the sentential complement. This is the mechanism by which TAGs are able to maintain long-distance dependencies (see Chapter 13), even over multiple embeddings (e.g. *Who did Bill expect to eat beans?*) or *who did Bill expect Mary to like?* See section 8.6.1 for details on this family. 20 verbs select this tree family.

Examples: *expect, see*

Van expects Bob to talk . Bob sees the harmonica fall .

Declarative tree: See Figure 6.45.

Figure 6.45: ECM Tree: $\beta Xnx0Vs1$

Other available trees: wh-moved subject, subject relative clause with and without comp, adjunct (gap-less) relative clause with and without comp/with PP pied-piping, imperative, NP gerund.

6.45 Idiom with V, D, and N anchors: Tnx0VDN1

Description: This tree family is selected by idiomatic phrases in which the verb, determiner, and NP are all frozen (as in *He kicked the bucket.*). Only a limited number of transformations are allowed, as compared to the normal transitive tree family (see section 6.2). Other idioms that have the same structure as *kick the bucket*, and that are limited to the same transformations would select this tree, while different tree families are used to handle other idioms. Note that *John kicked the bucket* is actually ambiguous, and would result in two parses - an idiomatic one (meaning that John died), and a compositional transitive one (meaning that there is an physical bucket that John hit with his foot). 1 idiom selects this family.

Examples: *kick the bucket*
Nixon kicked the bucket .

Declarative tree: See Figure 6.46.

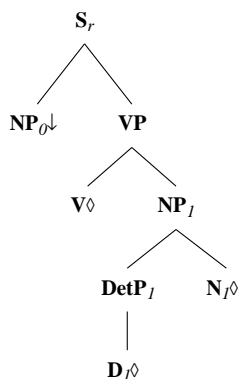


Figure 6.46: Declarative Transitive Idiom Tree: Tnx0VDN1

Other available trees: subject relative clause with and without comp, declarative, wh-moved subject, imperative, NP gerund, adjunct gapless relative with comp/with PP pied-piping, passive, w/wo by-phrase, wh-moved object of by-phrase, wh-moved by-phrase, relative (with and without comp) on subject of passive, PP relative.

6.46 Idiom with V, D, A, and N anchors: Tnx0VDAN1

Description: This tree family is selected by transitive idioms that are anchored by a verb, determiner, adjective, and noun. 19 idioms select this family.

Examples: *have a green thumb, sing a different tune*

Martha might have a green thumb, but it's uncertain after the death of all the plants.

After his conversion John sang a different tune.

Declarative tree: See Figure 6.47.

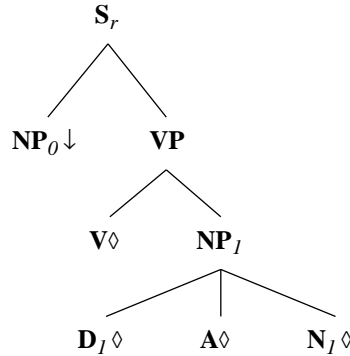


Figure 6.47: Declarative Idiom with V, D, A, and N Anchors Tree: $\alpha n x 0 V D A N 1$

Other available trees: Subject relative clause with and without comp, adjunct relative clause with comp/with PP pied-piping, wh-moved subject, imperative, NP gerund, passive without *by* phrase, passive with *by* phrase, passive with wh-moved object of *by* phrase, passive with wh-moved *by phrase*, passive with relative on object of *by* phrase with and without comp.

6.47 Idiom with V and N anchors: $T n x 0 V N 1$

Description: This tree family is selected by transitive idioms that are anchored by a verb and noun. 15 idioms select this family.

Examples: *draw blood, cry wolf*

Graham's retort drew blood.

The neglected boy cried wolf.

Declarative tree: See Figure 6.48.

Other available trees: Subject relative clause with and without comp, adjunct relative clause with comp/with PP pied-piping, wh-moved subject, imperative, NP gerund, passive without *by* phrase, passive with *by* phrase, passive with wh-moved object of *by* phrase, passive with wh-moved *by phrase*, passive with relative on object of *by* phrase with and without comp.

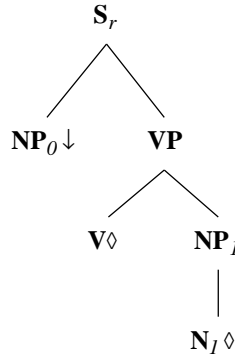


Figure 6.48: Declarative Idiom with V and N Anchors Tree: $\alpha nx0VN1$

6.48 Idiom with V, A, and N anchors: $Tnx0VAN1$

Description: This tree family is selected by transitive idioms that are anchored by a verb, adjective, and noun. 4 idioms select this family.

Examples: *break new ground, cry bloody murder*
The avant-garde film breaks new ground.
The investors cried bloody murder after the suspicious takeover.

Declarative tree: See Figure 6.49.

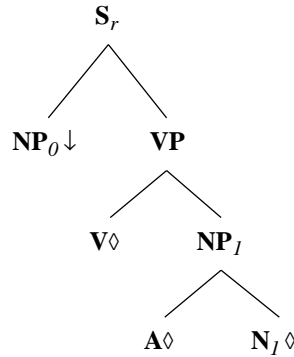


Figure 6.49: Declarative Idiom with V, A, and N Anchors Tree: $\alpha nx0VAN1$

Other available trees: Subject relative clause with and without comp, adjunct relative clause with comp/with PP pied-piping, wh-moved subject, imperative, NP gerund, passive without *by* phrase, passive with *by* phrase, passive with wh-moved object of *by* phrase, passive with wh-moved *by phrase*, passive with relative on object of *by* phrase with and without comp.

6.49 Idiom with V, D, A, N, and Prep anchors: Tnx0VDAN1Pnx2

Description: This tree family is selected by transitive idioms that are anchored by a verb, determiner, adjective, noun, and preposition. 6 idioms select this family.

Examples: *make a big deal about, make a great show of*
John made a big deal about a miniscule dent in his car.
The company made a big show of paying generous dividends.

Declarative tree: See Figure 6.50.

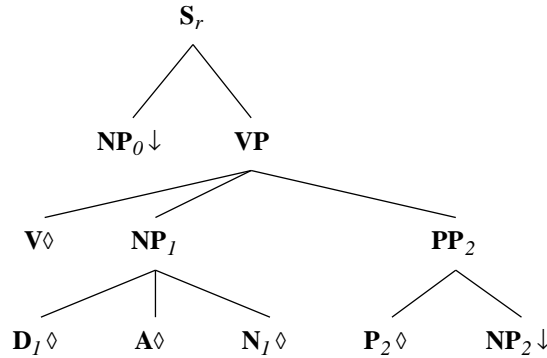


Figure 6.50: Declarative Idiom with V, D, A, N, and Prep Anchors Tree: α nx0VDAN1Pnx2

Other available trees: Subject relative clause with and without comp, adjunct relative clause with comp/with PP pied-piping, wh-moved subject, imperative, NP gerund, passive without *by* phrase, passive with *by* phrase, passive with wh-moved object of *by* phrase, passive with wh-moved *by* phrase, outer passive without *by* phrase, outer passive with *by* phrase, outer passive with wh-moved *by* phrase, outer passive with wh-moved object of *by* phrase, outer passive without *by* phrase with relative on the subject with and without comp, outer passive with *by* phrase with relative on subject with and without comp.

6.50 Idiom with V, A, N, and Prep anchors: Tnx0VAN1Pnx2

Description: This tree family is selected by transitive idioms that are anchored by a verb, adjective, noun, and preposition. 3 idioms select this family.

Examples: *make short work of*
John made short work of the glazed ham.

Declarative tree: See Figure 6.51.

Other available trees: Subject relative clause with and without comp, adjunct relative clause with comp/with PP pied-piping, wh-moved subject, imperative, NP gerund, passive without *by* phrase, passive with *by* phrase, passive with wh-moved object of *by* phrase, passive

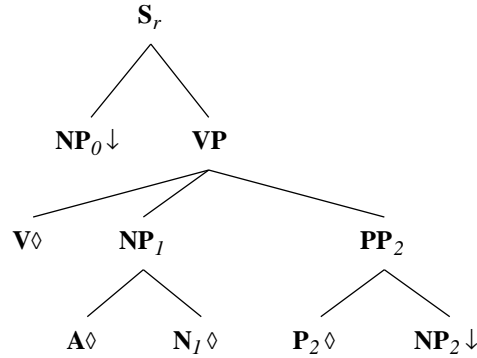


Figure 6.51: Declarative Idiom with V, A, N, and Prep Anchors Tree: $\alpha n x 0 V A N 1 P n x 2$

with wh-moved *by phrase*, outer passive without *by* phrase, outer passive with *by* phrase, outer passive with wh-moved *by* phrase, outer passive with wh-moved object of *by* phrase, outer passive without *by* phrase with relative on the subject with and without comp, outer passive with *by* phrase with relative on subject with and without comp.

6.51 Idiom with V, N, and Prep anchors: $T n x 0 V N 1 P n x 2$

Description: This tree family is selected by transitive idioms that are anchored by a verb, noun, and preposition. 6 idioms select this family.

Examples: *look daggers at*, *keep track of*
Maria looked daggers at her ex-husband across the courtroom.
The company kept track of its inventory.

Declarative tree: See Figure 6.52.

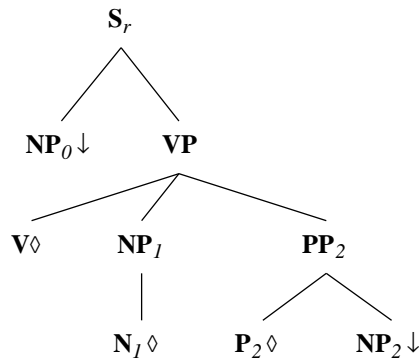


Figure 6.52: Declarative Idiom with V, N, and Prep Anchors Tree: $\alpha n x 0 V N 1 P n x 2$

Other available trees: Subject relative clause with and without comp, adjunct relative clause with comp/with PP pied-piping, wh-moved subject, imperative, NP gerund, passive with-

out *by* phrase, passive with *by* phrase, passive with wh-moved object of *by* phrase, passive with wh-moved *by phrase*, outer passive without *by* phrase, outer passive with *by* phrase, outer passive with wh-moved *by phrase*, outer passive with wh-moved object of *by* phrase, outer passive without *by* phrase with relative on the subject with and without comp, outer passive with *by* phrase with relative on subject with and without comp.

6.52 Idiom with V, D, N, and Prep anchors: Tnx0VDN1Pnx2

Description: This tree family is selected by transitive idioms that are anchored by a verb, determiner, noun, and preposition. 17 idioms select this family.

Examples: *make a mess of, keep the lid on*
John made a mess out of his new suit.
The tabloid didn't keep a lid on the imminent celebrity nuptials.

Declarative tree: See Figure 6.53.

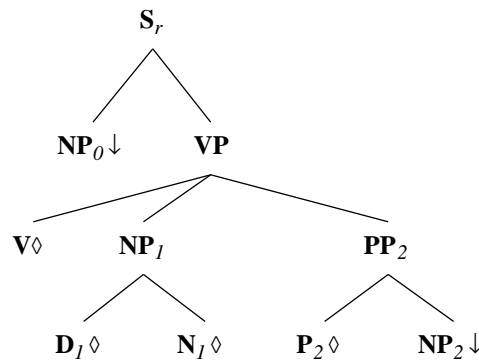


Figure 6.53: Declarative Idiom with V, D, N, and Prep Anchors Tree: αnx0VDN1Pnx2

Other available trees: Subject relative clause with and without comp, adjunct relative clause with comp/with PP pied-piping, wh-moved subject, imperative, NP gerund, passive without *by* phrase, passive with *by* phrase, passive with wh-moved object of *by* phrase, passive with wh-moved *by phrase*, outer passive without *by* phrase, outer passive with *by* phrase, outer passive with wh-moved *by phrase*, outer passive with wh-moved object of *by* phrase, outer passive without *by* phrase with relative on the subject with and without comp, outer passive with *by* phrase with relative on subject with and without comp.

Chapter 7

Ergatives

Verbs in English that are termed ergative display the kind of alternation shown in the sentences in (7) below.

- (7) The sun melted the ice .
The ice melted .

The pattern of ergative pairs as seen in (7) is for the object of the transitive sentence to be the subject of the intransitive sentence. The literature discussing such pairs is based largely on syntactic models that involve movement, particularly GB. Within that framework two basic approaches are discussed:

- **Derived Intransitive**

The intransitive member of the ergative pair is derived through processes of movement and deletion from:

- a transitive D-structure [Burzio, 1986]; or
- transitive lexical structure [Hale and Keyser, 1986; Hale and Keyser, 1987]

- **Pure Intransitive**

The intransitive member is intransitive at all levels of the syntax and the lexicon and is not related to the transitive member syntactically or lexically [Napoli, 1988].

The Derived Intransitive approach's notions of movement in the lexicon or in the grammar are not represented as such in the XTAG grammar. However, distinctions drawn in these arguments can be translated to the FB-LTAG framework. In the XTAG grammar the difference between these two approaches is not a matter of movement but rather a question of tree family membership. The relation between sentences represented in terms of movement in other frameworks is represented in XTAG by membership in the same tree family. Wh-questions and their indicative counterparts are one example of this. Adopting the Pure Intransitive approach suggested by [Napoli, 1988] would mean placing the intransitive ergatives in a tree family with other intransitive verbs and separate from the transitive variants of the same verbs. This would result in a grammar that represented intransitive ergatives as more closely related to other intransitives than to their transitive counterparts. The only hint of the relation between the

intransitive ergatives and the transitive ergatives would be that ergative verbs would select both tree families. While this is a workable solution, it is an unattractive one for the English XTAG grammar because semantic coherence is implicitly associated with tree families in our analysis of other constructions. In particular, constancy in thematic role is represented by constancy in node names across sentence types within a tree family. For example, if the object of a declarative tree is NP₁ the subject of the passive tree(s) in that family will also be NP₁.

The analysis that has been implemented in the English XTAG grammar is an adaptation of the Derived Intransitive approach. The ergative verbs select one family, Tnx0Vnx1, that contains both transitive and intransitive trees. The <trans> feature appears on the intransitive ergative trees with the value – and on the transitive trees with the value +. This creates the two possibilities needed to account for the data.

- **intransitive ergative/transitive alternation.** These verbs have transitive and intransitive variants as shown in sentences (8) and (9).

(8) The sun melted the ice cream .

(9) The ice cream melted .

In the English XTAG grammar, verbs with this behavior are left unspecified as to value for the <trans> feature. This lack of specification allows these verbs to anchor either type of tree in the Tnx0Vnx1 tree family because the unspecified <trans> value of the verb can unify with either + or – values in the trees.

- **transitive only.** Verbs of this type select only the transitive trees and do not allow intransitive ergative variants as in the pattern show in sentences (10) and (11).

(10) Elmo borrowed a book .

(11) *A book borrowed .

The restriction to selecting only transitive trees is accomplished by setting the <trans> feature value to + for these verbs.

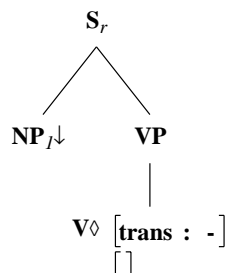


Figure 7.1: Ergative Tree: α Enx1V

The declarative ergative tree is shown in Figure 7.1 with the <trans> feature displayed. Note that the index of the subject NP indicates that it originated as the object of the verb.

Chapter 8

Sentential Subjects and Sentential Complements

In the XTAG grammar, arguments of a lexical item, including subjects, appear in the initial tree anchored by that lexical item. A sentential argument appears as an S node in the appropriate position within an elementary tree anchored by the lexical item that selects it. This is the case for sentential complements of verbs, prepositions and nouns and for sentential subjects. The distribution of complementizers in English is intertwined with the distribution of embedded sentences. A successful analysis of complementizers in English must handle both the cooccurrence restrictions between complementizers and various types of clauses, and the distribution of the clauses themselves, in both subject and complement positions.

8.1 S or VP complements?

Two comparable grammatical formalisms, Generalized Phrase Structure Grammar (GPSG) [Gazdar *et al.*, 1985] and Head-driven Phrase Structure Grammar (HPSG) [Pollard and Sag, 1994], have rather different treatments of sentential complements (S-comps). They both treat embedded sentences as VP's with subjects, which generates the correct structures but misses the generalization that S's behave similarly in both matrix and embedded environments, and VP's behave quite differently. Neither account has PRO subjects of infinitival clauses— they have subjectless VP's instead. GPSG has a complete complementizer system, which appears to cover the same range of data as our analysis. It is not clear what sort of complementizer analysis could be implemented in HPSG.

Following standard GB approach, the English XTAG grammar does not allow VP complements but treats verb-anchored structures without overt subjects as having PRO subjects. Thus, indicative clauses, infinitives and gerunds all have a uniform treatment as embedded clauses using the same trees under this approach. Furthermore, our analysis is able to preserve the selectional and distributional distinction between S's and VP's, in the spirit of GB theories, without having to posit 'extra' empty categories.¹ Consider the alternation between *that* and the null complementizer², shown in sentences (12) and (13).

¹i.e. empty complementizers. We do have PRO and NP traces in the grammar.

²Although we will continue to refer to 'null' complementizers, in our analysis this is actually the absence of

(12) He hopes \emptyset Muriel wins.

(13) He hopes that Muriel wins.

In GB both *Muriel wins* in (12) and *that Muriel wins* in (13) are CPs even though there is no overt complementizer to head the phrase in (12). Our grammar does not distinguish by category label between the phrases that would be labeled in GB as IP and CP. We label both of these phrases S. The difference between these two levels is the presence or absence of the complementizer (or extracted WH constituent), and is represented in our system as a difference in feature values (here, of the **<comp>** feature), and the presence of the additional structure contributed by the complementizer or extracted constituent. This illustrates an important distinction in XTAG, that between features and node labels. Because we have a sophisticated feature system, we are able to make fine-grained distinctions between nodes with the same label which in another system might have to be realized by using distinguishing node labels.

8.2 Complementizers and Embedded Clauses in English: The Data

Verbs selecting sentential complements (or subjects) place restrictions on their complements, in particular, on the form of the embedded verb phrase.³ Furthermore, complementizers are constrained to appear with certain types of clauses, again, based primarily on the form of the embedded VP. For example, *hope* selects both indicative and infinitival complements. With an indicative complement, it may only have *that* or null as possible complementizers; with an infinitival complement, it may only have a null complementizer. Verbs that allow wh+ complementizers, such as *ask*, can take *whether* and *if* as complementizers. The possible combinations of complementizers and clause types is summarized in Table 8.1.

As can be seen in Table 8.1, sentential subjects differ from sentential complements in requiring the complementizer *that* for all indicative and subjunctive clauses. In sentential complements, *that* often varies freely with a null complementizer, as illustrated in (14)-(19).

(14) Christy hopes that Mike wins.

(15) Christy hopes Mike wins.

(16) Dania thinks that Newt is a liar.

(17) Dania thinks Newt is a liar.

(18) That Helms won so easily annoyed me.

(19) *Helms won so easily annoyed me.

a complementizer.

³Other considerations, such as the relationship between the tense/aspect of the matrix clause and the tense/aspect of a complement clause are also important but are not currently addressed in the current English XTAG grammar.

Complementizer:		that	whether	if	for	null
Clause type						
indicative	subject	Yes	Yes	No	No	No
	complement	Yes	Yes	Yes	No	Yes
infinitive	subject	No	Yes	No	Yes	Yes
	complement	No	Yes	No	Yes	Yes
subjunctive	subject	Yes	No	No	No	No
	complement	Yes	No	No	No	Yes
gerundive ⁴	complement	No	No	No	No	Yes
base	complement	No	No	No	No	Yes
small clause	complement	No	No	No	No	Yes

Table 8.1: Summary of Complementizer and Clause Combinations

Another fact which must be accounted for in the analysis is that in infinitival clauses, the complementizer *for* must appear with an overt subject NP, whereas a complementizer-less infinitival clause never has an overt subject, as shown in (20)-(23). (See section 8.5 for more discussion of the case assignment issues relating to this construction.)

(20) To lose would be awful.

(21) For Penn to lose would be awful.

(22) *For to lose would be awful.

(23) *Penn to lose would be awful.

In addition, some verbs select $\langle \mathbf{wh} \rangle = +$ complements (either questions or clauses with *whether* or *if*) [Grimshaw, 1990]:

(24) Jesse wondered who left.

(25) Jesse wondered if Barry left.

(26) Jesse wondered whether to leave.

(27) Jesse wondered whether Barry left.

(28) *Jesse thought who left.

(29) *Jesse thought if Barry left.

(30) *Jesse thought whether to leave.

(31) *Jesse thought whether Barry left.

⁴Most gerundive phrases are treated as NP's. In fact, all gerundive subjects are treated as NP's, and the only gerundive complements which receive a sentential parse are those for which there is no corresponding NP parse. This was done to reduce duplication of parses. See Chapter 17 for further discussion of gerunds.

8.3 Features Required

As we have seen above, clauses may be $\langle \mathbf{wh} \rangle = +$ or $\langle \mathbf{wh} \rangle = -$, may have one of several complementizers or no complementizer, and can be of various clause types. The XTAG analysis uses three features to capture these possibilities: $\langle \mathbf{comp} \rangle$ for the variation in complementizers, $\langle \mathbf{wh} \rangle$ for the question vs. non-question alternation and $\langle \mathbf{mode} \rangle$ ⁵ for clause types. In addition to these three features, the $\langle \mathbf{assign-comp} \rangle$ feature represents complementizer requirements of the embedded verb. More detailed discussion of the $\langle \mathbf{assign-comp} \rangle$ feature appears below in the discussions of sentential subjects and of infinitives. The four features and their possible values are shown in Table 8.2.

Feature	Values
$\langle \mathbf{comp} \rangle$	that, if, whether, for, rel, nil
$\langle \mathbf{mode} \rangle$	ind, inf, subjnt, ger, base, ppart, nom/prep
$\langle \mathbf{assign-comp} \rangle$	that, if, whether, for, rel, ind_nil, inf_nil
$\langle \mathbf{wh} \rangle$	+,-

Table 8.2: Summary of Relevant Features

8.4 Distribution of Complementizers

Like other non-arguments, complementizers anchor an auxiliary tree (shown in Figure 8.1) and adjoin to elementary clausal trees. The auxiliary tree for complementizers is the only alternative to having a complementizer position ‘built into’ every sentential tree. The latter choice would mean having an empty complementizer substitute into every matrix sentence and a complementizerless embedded sentence to fill the substitution node. Our choice follows the XTAG principle that initial trees consist only of the arguments of the anchor⁶ – the S tree does not contain a slot for a complementizer, and the β COMP tree has only one argument, an S with particular features determined by the complementizer. Complementizers select the type of clause to which they adjoin through constraints on the $\langle \mathbf{mode} \rangle$ feature of the S foot node in the tree shown in Figure 8.1. These features also pass up to the root node, so that they are ‘visible’ to the tree where the embedded sentence adjoins/substitutes.

The grammar handles the following complementizers: *that*, *whether*, *if*, *for*, and no complementizer, and the clause types: indicative, infinitival, gerundive, past participial, subjunctive and small clause (**nom/prep**). The $\langle \mathbf{comp} \rangle$ feature in a clausal tree reflects the value of the complementizer if one has adjoined to the clause.

The $\langle \mathbf{comp} \rangle$ and $\langle \mathbf{wh} \rangle$ features receive their root node values from the particular complementizer which anchors the tree. The β COMPs tree adjoins to an S node with the feature $\langle \mathbf{comp} \rangle = \mathbf{nil}$; this feature indicates that the tree does not already **have** a complementizer adjoined to it.⁷ We ensure that there are no stacked complementizers by requiring the foot node of β COMPs to have $\langle \mathbf{comp} \rangle = \mathbf{nil}$.

⁵ $\langle \mathbf{mode} \rangle$ actually conflates several types of information, in particular verb form and mood.

⁶See section 4.2 for a discussion of the difference between complements and adjuncts in the XTAG grammar.

⁷Because root S’s cannot have complementizers, the parser checks that the root S has $\langle \mathbf{comp} \rangle = \mathbf{nil}$ at the end of the derivation, when the S is also checked for a tensed verb.

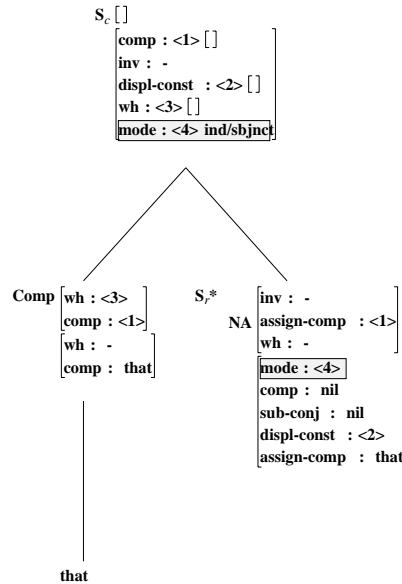


Figure 8.1: Tree β COMPs, anchored by *that*

8.5 Case assignment, *for* and the two *to*'s

The **<assign-comp>** feature is used to represent the requirements of particular types of clauses for particular complementizers. So while the **<comp>** feature represents constraints originating from the VP dominating the clause, the **<assign-comp>** feature represents constraints originating from the highest VP in the clause. **<assign-comp>** is used to control the appearance of subjects in infinitival clauses (see discussion of ECM constructions in 8.6.1), to block bare indicative sentential subjects (bare infinitival subjects are allowed), and to block ‘that-trace’ violations.

Examples (33), (34) and (35) show that an accusative case subject is obligatory in an infinitive clause if the complementizer *for* is present. The infinitive clauses in (32) is analyzed in the English XTAG grammar as having a PRO subject.

(32) Christy wants to pass the exam.

(33) Mike wants for her to pass the exam.

(34) *Mike wants for she to pass the exam.

(35) *Christy wants for to pass the exam.

The *for-to* construction is particularly illustrative of the difficulties and benefits faced in using a lexicalized grammar. It is commonly accepted that *for* behaves as a case-assigning complementizer in this construction, assigning accusative case to the ‘subject’ of the clause since the infinitival verb does not assign case to its subject position. However, in our featurized grammar, the absence of a feature licenses anything, so we must have overt null case assigned

by infinitives to ensure the correct distribution of PRO subjects. (See section 4.4 for more discussion of case assignment.) This null case assignment clashes with accusative case assignment if we simply add *for* as a standard complementizer, since NP's (including PRO) are drawn from the lexicon already marked for case. Thus, we must use the **<assign-comp>** feature to pass information about the verb up to the root of the embedded sentence. To capture these facts, two infinitive *to*'s are posited. One infinitive *to* has **<assign-case>=none** which forces a PRO subject, and **<assign-comp>=inf_nil** which prevents *for* from adjoining. The other infinitive *to* has no value at all for **<assign-case>** and has **<assign-comp>=for/ecm**, so that it can only occur either with the complementizer *for* or with ECM constructions. In those instances either *for* or the ECM verb supplies the **<assign-case>** value, assigning accusative case to the overt subject.

8.6 Sentential Complements of Verbs

Tree families: Tnx0Vs1, Tnx0Vnx1s2, TItVnx1s2, TItVpnx1s2, TItVad1s2.

Verbs that select sentential complements restrict the **<mode>** and **<comp>** values for those complements. Since with very few exceptions⁸ long distance extraction is possible from sentential complements, the S complement nodes are adjunction nodes. Figure 8.2 shows the declarative tree for sentential complements, anchored by *think*.

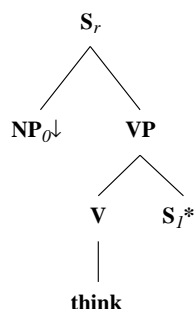


Figure 8.2: Sentential complement tree: β_{nx0Vs1}

The need for an adjunction node rather than a substitution node at S_1 may not be obvious until one considers the derivation of sentences with long distance extractions. For example, the declarative in (36) is derived by adjoining the tree in Figure 8.3(b) to the S_1 node of the tree in Figure 8.3(a). Since there are no bottom features on S_1 , the same final result could have been achieved with a substitution node at S_1 .

(36) The emu thinks that the aardvark smells terrible.

However, adjunction is crucial in deriving sentences with long distance extraction, as in sentences (37) and (38).

(37) Who does the emu think smells terrible?

⁸For example, long distance extraction is not possible from the S complement in it-clefts.

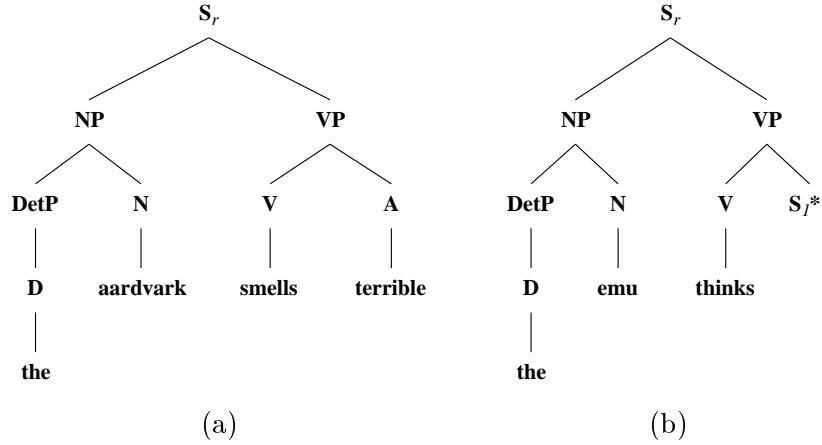


Figure 8.3: Trees for *The emu thinks that the aardvark smells terrible*.

(38) Who did the elephant think the panda heard the emu say smells terrible?

The example in (37) is derived from the trees for *who smells terrible?* shown in Figure 8.4 and *the emu thinks* S shown in Figure 8.3(b), by adjoining the latter at the S_r node of the former.⁹ This process is recursive, allowing sentences like (38). Such a representation has been shown by [Kroch and Joshi, 1985] to be well-suited for describing unbounded dependencies.

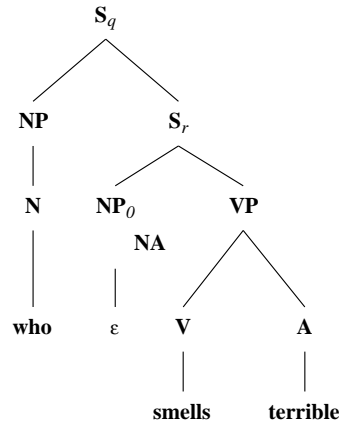


Figure 8.4: Tree for *Who smells terrible?*

In English, a complementizer may not appear on a complement with an extracted subject (the ‘that-trace’ configuration). This phenomenon is illustrated in (39)-(41):

(39) Which animal did the giraffe say that he likes?

(40) *Which animal did the giraffe say that likes him?

(41) Which animal did the giraffe say likes him?

⁹See Chapter 20 for a discussion of do-support.

These sentences are derived in XTAG by adjoining the tree for *did the giraffe say S* at the S_r node of the tree for either *which animal likes him* (to yield sentence (41)) or *which animal he likes* (to yield sentence (39)). That-trace violations are blocked by the presence of the feature $\langle \text{assign-comp} \rangle = \text{inf_nil/ind_nil/ecm}$ feature on the bottom of the S_r node of trees with extracted subjects (W0), i.e. those used in sentences such as (40) and (41). If a complementizer tree, β COMPs, adjoins to a subject extraction tree at S_r , its $\langle \text{assign-comp} \rangle = \text{that/whether/for/if}$ feature will clash and the derivation will fail. If there is no complementizer, there is no feature clash, and this will permit the derivation of sentences like (41), or of ECM constructions, in which case the ECM verb will have $\langle \text{assign-comp} \rangle = \text{ecm}$ (see section 8.6.1 for more discussion of the ECM case). Complementizers may adjoin normally to object extraction trees such as those used in sentence (39), and so object extraction trees have no value for the $\langle \text{assign-comp} \rangle$ feature.

In the case of indirect questions, subadjacency follows from the principle that a given tree cannot contain more than one wh-element. Extraction out of an indirect question is ruled out because a sentence like:

(42) * Who_i do you wonder who_j e_j loves e_i ?

would have to be derived from the adjunction of *do you wonder* into *who_i who_j e_j loves e_i*, which is an ill-formed elementary tree.¹⁰

8.6.1 Exceptional Case Marking Verbs

Tree family: TXnx0Vs1 Exceptional Case Marking verbs are those which assign accusative case to the subject of the sentential complement. This is in contrast to verbs in the Tnx0Vnx1s2 family (section 6.6), which assign accusative case to an NP which is not part of the sentential complement.

The subject of an ECM infinitive complement is assigned accusative case in a manner analogous to that of a subject in a *for-to* construction, as described in section 8.5. As in the *for-to* case, the ECM verb assigns accusative case into the subject of the lower infinitive, and so the infinitive uses the *to* which has no value for $\langle \text{assign-case} \rangle$ and has $\langle \text{assign-comp} \rangle = \text{for/ecm}$. The ECM verb has $\langle \text{assign-comp} \rangle = \text{ecm}$ and $\langle \text{assign-case} \rangle = \text{acc}$ on its foot. The former allows the $\langle \text{assign-comp} \rangle$ features of the ECM verb and the *to* tree to unify, and so be used together, and the latter assigns the accusative case to the lower subject.

Figure 8.5 shows the declarative tree for the tree for the TXnx0Vs1 family, in this case anchored by *expects*. Figure 8.6 shows a parse for *Van expects Bob to talk*

The ECM and *for-to* cases are analogous in how they are used together with the correct infinitival *to* to assign accusative case to the subject of the lower infinitive. However, they are different in that *for* is blocked along with other complementizers in subject extraction contexts, as discussed in section 8.6, as in (43), while subject extraction is compatible with ECM cases, as in (44).

(43) *What child did the giraffe ask for to leave?

¹⁰This does not mean that elementary trees with more than one gap should be ruled out across the grammar. Such trees might be required for dealing with parasitic gaps or gaps in coordinated structures.

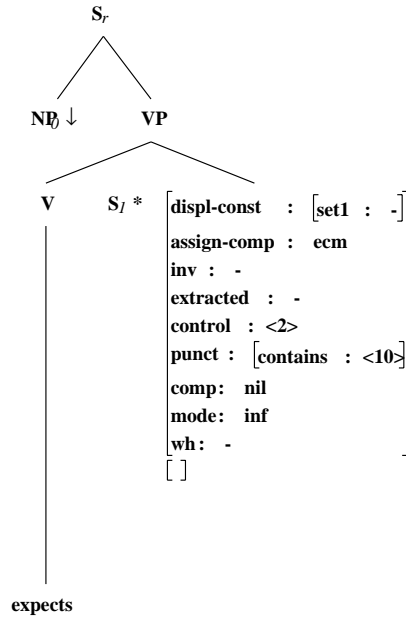


Figure 8.5: ECM tree: βX_{nx0Vs1}

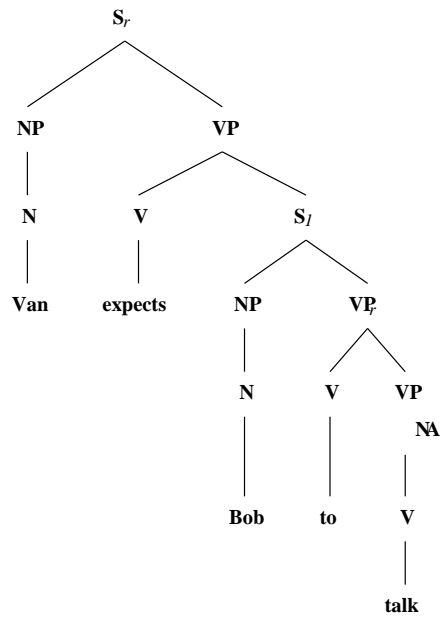


Figure 8.6: Sample ECM parse

(44) Who did Bill expect to eat beans?

Sentence (43) is ruled out by the $\langle \text{assign-comp} \rangle = \text{inf_nil/ind_nil/ecm}$ feature on the subject extraction tree for *ask*, since the $\langle \text{assign-comp} \rangle = \text{for}$ feature from the *for* tree will

fail to unify. However, (44) will be allowed since $\langle \text{assign-comp} \rangle = \text{ecm}$ feature on the *expect* tree will unify with the foot of the ECM verb tree. The use of features allows the ECM and *for-to* constructions to act the same for exceptional case assignment, while also being distinguished for that-trace violations.

Verbs that take bare infinitives, as in (45), are also treated as ECM verbs, the only difference being that their foot feature has $\langle \text{mode} \rangle = \text{base}$ instead of $\langle \text{mode} \rangle = \text{inf}$. Since the complement does not have *to*, there is no question of using the *to* tree for allowing accusative case to be assigned. Instead, verbs with $\langle \text{mode} \rangle = \text{base}$ allow either accusative or nominative case to be assigned to the subject, and the foot of the ECM bare infinitive tree forces it to be accusative by its $\langle \text{assign-case} \rangle = \text{acc}$ value at its foot node unifies with the $\langle \text{assign-case} \rangle = \text{nom/acc}$ value of the bare infinitive clause.

(45) Bob sees the harmonica fall.

The trees in the TXnx0Vs1 family are generally parallel to those in the Tnx0Vs1 family, except for the $\langle \text{assign-case} \rangle$ and $\langle \text{assign-comp} \rangle$ values on the foot nodes. However, the TXnx0Vs1 family also includes a tree for the passive, which of course is not included in the Tnx0Vs1 family. Unlike all the other trees in the TXnx0Vs1 family, the passive tree is not rooted in S, and is instead a VP auxiliary tree. Since the subject of the infinitive is not thematically selected by the ECM verb, it is not part of the ECM verb's tree, and so it cannot be part of the passive tree. Therefore, the passive acts as a raising verb (see section 9.3). For example, to derive (47), the tree in Figure 8.7 would adjoin into a derivation for *Bob to talk* at the VP node (and the $\langle \text{mode} \rangle = \text{passive}$ feature, not shown, forces the auxiliary to adjoin in, as for other passives, as described in chapter 12).

(46) Van expects Bob to talk.

(47) Bob was expected to talk.

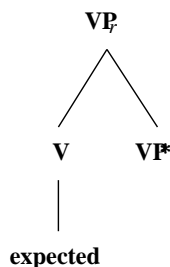


Figure 8.7: ECM passive

It has long been noted that passives of both full and bare infinitive ECM constructions are full infinitives, as in (47) and (49).

(48) Bob sees the harmonica fall.

(49) The harmonica was seen to fall.

(50) *The harmonica was seen fall.

Under the TAG ECM analysis, this fact is easy to implement. The foot node of the ECM passive tree is simply set to have $\langle \mathbf{mode} \rangle = \mathbf{inf}$, which prevents the derivation of (50). Therefore, for all the other trees in the family, foot nodes are set to have $\langle \mathbf{mode} \rangle = \mathbf{base}$ or $\langle \mathbf{mode} \rangle = \mathbf{inf}$ depending on whether it is a bare infinitive or not. These foot nodes are all S nodes. The VP foot node of the passive tree, however, has $\langle \mathbf{mode} \rangle = \mathbf{inf}$ regardless.

8.7 Sentential Subjects

Tree families: Ts0Vnx1, Ts0Ax1, Ts0N1, Ts0PNx1, Ts0ARBPNx1, Ts0PPnx1, Ts0PNaPNx1, Ts0V, Ts0Vtonx1, Ts0NPnx1, Ts0APnx1, Ts0A1s1.

Verbs that select sentential subjects anchor trees that have an S node in the subject position rather than an NP node. Since extraction is not possible from sentential subjects, they are implemented as substitution nodes in the English XTAG grammar. Restrictions on sentential subjects, such as the required *that* complementizer for indicatives, are enforced by feature values specified on the S substitution node in the elementary tree.

Sentential subjects behave essentially like sentential complements, with a few exceptions. In general, all verbs which license sentential subjects license the same set of clause types. Thus, unlike sentential complement verbs which select particular complementizers and clause types, the matrix verbs licensing sentential subjects merely license the S argument. Information about the complementizer or embedded verb is located in the tree features, rather than in the features of each verb selecting that tree. Thus, all sentential subject trees have the same $\langle \mathbf{mode} \rangle$, $\langle \mathbf{comp} \rangle$ and $\langle \mathbf{assign-comp} \rangle$ values shown in Figure 8.8(a).

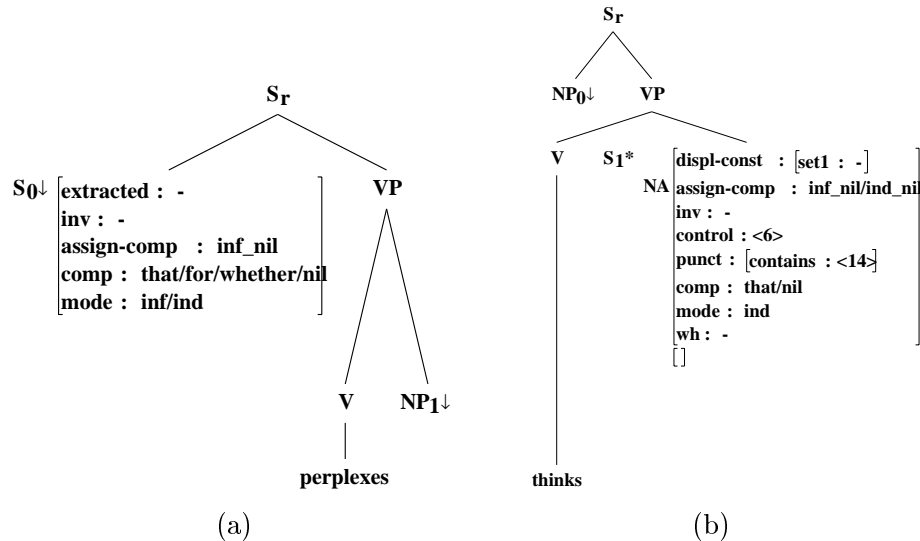


Figure 8.8: Comparison of $\langle \mathbf{assign-comp} \rangle$ values for sentential subjects: $\alpha s0Vnx1$ (a) and sentential complements: $\beta nx0Vs1$ (b)

The major difference in clause types licensed by S-subjs and S-comps is that indicative S-subjs obligatorily have a complementizer (see examples in section 8.2). The $\langle \mathbf{assign-comp} \rangle$

feature is used here to license a null complementizer for infinitival but not indicative clauses. $\langle \mathbf{assign-comp} \rangle$ has the same possible values as $\langle \mathbf{comp} \rangle$, with the exception that the **nil** value is ‘split’ into **ind_nil** and **inf_nil**. This difference in feature values is illustrated in Figure 8.8.

Another minor difference is that *whether* but not *if* is grammatical with S-subjs.¹¹ Thus, *if* is not among the $\langle \mathbf{comp} \rangle$ values allowed in S-subjs. The final difference from S-comps is that there are no S-subjs with $\langle \mathbf{mode} \rangle = \mathbf{ger}$. As noted in footnote 4 of this chapter, gerundive complements are only allowed when there is no corresponding NP parse. In the case of gerundive S-subjs, there is always an NP parse available.

8.8 Nouns and Prepositions taking Sentential Complements

Trees: αNXNs , βvxPs , βPss , βnxPs , Tnx0N1s1 , Tnx0A1s1 .

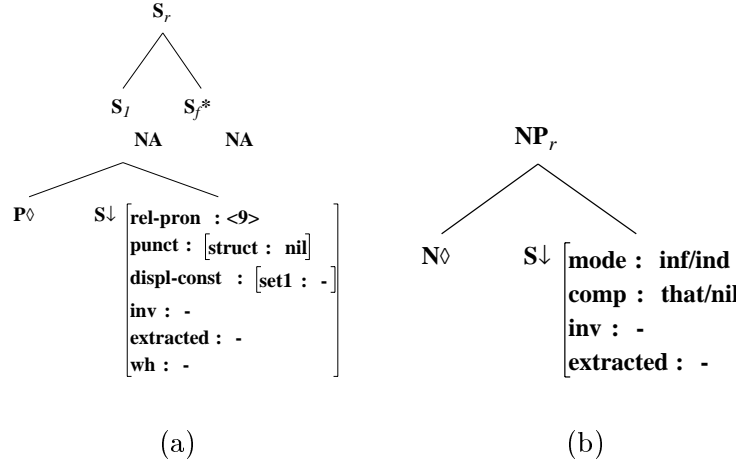


Figure 8.9: Sample trees for preposition: βPss (a) and noun: αNXNs (b) taking sentential complements

Prepositions and nouns can also select sentential complements, using the trees listed above. These trees use the $\langle \mathbf{mode} \rangle$ and $\langle \mathbf{comp} \rangle$ features as shown in Figure 8.9. For example, the noun *claim* takes only indicative complements with *that*, while the preposition *with* takes small clause complements, as seen in sentences (51)-(54).

(51) Beth’s claim that Clove was a smart dog....

(52) *Beth’s claim that Clove a smart dog...

(53) Dania wasn’t getting any sleep with Doug sick.

(54) *Dania wasn’t getting any sleep with Doug was sick.

¹¹Some speakers also find *if* as a complementizer only marginally grammatical in S-comps.

8.9 PRO control

8.9.1 Types of control

In the literature on control, two types are often distinguished: obligatory control, as in sentences (55), (56), (57), and (58) and optional control, as in sentence (59).

(55) Srini_i promised Mickey_i [PRO_i to leave].

(56) Srini persuaded Mickey_i [PRO_i to leave].

(57) Srini_i wanted [PRO_i to leave].

(58) Christy_i left the party early [PRO_i to go to the airport].

(59) [PRO_{arb/i} to dance] is important for Bill_i.

At present, an analysis for obligatory control into complement clauses (as in sentences (55), (56), and (57)) has been implemented. An analysis for cases of obligatory control into adjunct clauses and optional control exists and can be found in [Bhatt, 1994].

8.9.2 A feature-based analysis of PRO control

The analysis for obligatory control involves co-indexation of the control feature of the NP anchored by PRO to the control feature of the controller. A feature equation in the tree anchored by the control verb co-indexes the control feature of the controlling NP with the foot node of the tree. All sentential trees have a co-indexed control feature from the root S to the subject NP.

When the tree containing the controller adjoins onto the complement clause tree containing the PRO, the features of the foot node of the auxiliary tree are unified with the bottom features of the root node of the complement clause tree containing the PRO. This leads to the control feature of the controller being co-indexed with the control feature of the PRO.

Depending on the choice of the controlling verb, the control propagation paths in the auxiliary trees are different. In the case of subject control (as in sentence (56)), the subject NP and the foot node are have co-indexed control features, while for object control (e.g. sentence (55)), the object NP and the foot node are co-indexed for control. Among verbs that belong to the Tnx0Vnx1s2 family, i.e. verbs that take an NP object and a clausal complement, subject-control verbs form a distinct minority, *promise* being the only commonly used verb in this class.

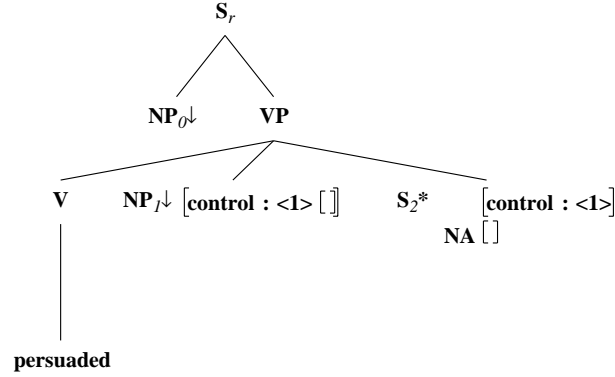
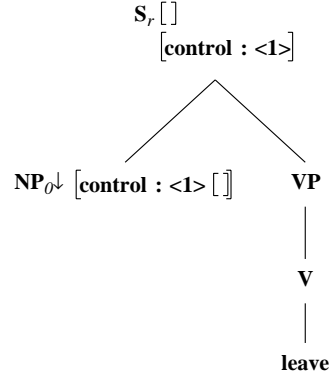
Consider the derivation of sentence (56). The auxiliary tree for *persuade*, shown in Figure 8.10, has the following feature equation (60).

$$(60) \text{NP}_1:\langle \mathbf{control} \rangle = \text{S}_2.\text{t}:\langle \mathbf{control} \rangle$$

The auxiliary tree adjoins into the tree for *leave*, shown in Figure 8.11, which has the following feature equation (61).

$$(61) \text{S}_r.\text{b}:\langle \mathbf{control} \rangle = \text{NP}_0.\text{t}:\langle \mathbf{control} \rangle$$

Since the adjunction takes place at the root node (S_r) of the *leave* tree, after unification, NP_1 of the *persuade* tree and NP_0 of the *leave* tree share a control feature. The resulting derived and derivation trees are shown in Figures 8.12 and 8.13.


 Figure 8.10: Tree for *persuaded*

 Figure 8.11: Tree for *leave*

8.9.3 The nature of the control feature

The control feature does not have any value and is used only for co-indexing purposes. If two NPs have their control features co-indexed, it means that they are participating in a relationship of control; the c-commanding NP controls the c-commanded NP.

8.9.4 Long-distance transmission of control features

Cases involving embedded infinitival complements with PRO subjects such as (62) can also be handled.

(62) John_{*i*} wants [PRO_{*i*} to want [PRO_{*i*} to dance]].

The control feature of ‘John’ and the two PRO’s all get co-indexed. This treatment might appear to lead to a problem. Consider (63):

(63) John_{**i*} wants [Mary_{*i*} to want [PRO_{*i*} to dance]].

If both the ‘want’ trees have the control feature of their subject co-indexed to their foot nodes, we would have a situation where the PRO is co-indexed for control feature with ‘John’,

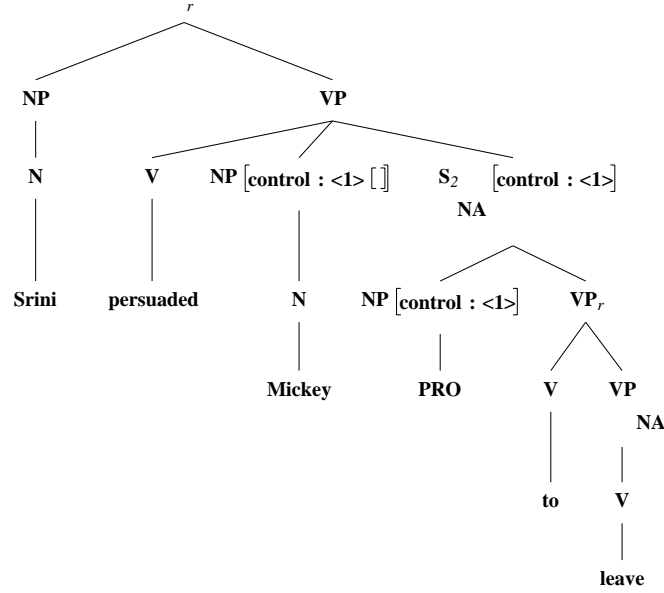


Figure 8.12: Derived tree for *Srini persuaded Mickey to leave*

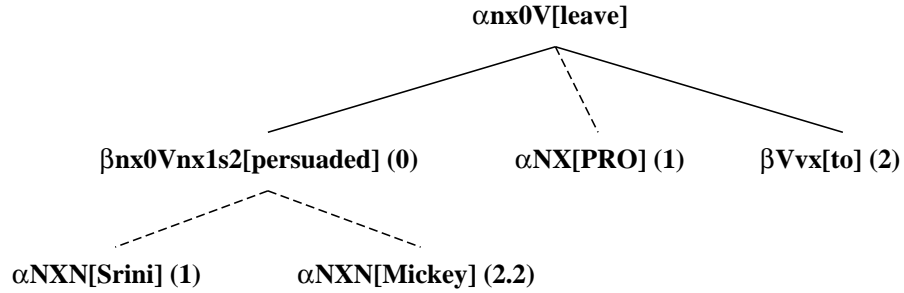


Figure 8.13: Derivation tree for *Srini persuaded Mickey to leave*

as well as with ‘Mary’. Note that the higher ‘want’ in (62) is *want_{ECM}* - it assigns case to the subject of the lower clause while the lower ‘want’ in (62) is not. Subject control is restricted to non-ECM (Exceptional Case Marking) verbs that take infinitival complements. Since the two ‘want’s in (62) are different with respect to their control (and other) properties, the control feature of PRO stops at ‘Mary’ and is not transmitted to the higher clause.

8.9.5 Locality constraints on control

PRO control obeys locality constraints. The controller for PRO has to be in the immediately higher clause. Consider the ungrammatical sentence (64) ((64) is ungrammatical only with the co-indexing indicated below).

(64) * John_i wants [PRO_i to persuade Mary_i [PRO_i to dance]]

However, such a derivation is ruled out automatically by the mechanisms of a TAG derivation and feature unification. Suppose it was possible to first compose the *want* tree with the *dance* tree and then insert the *persuade* tree. (This is not possible in the XTAG grammar because of the convention that auxiliary trees have NA (Null Adjunction) constraints on their foot nodes.) Even then, at the end of the derivation the control feature of the subject of *want* would end up co-indexed with the PRO subject of *persuade* and the control feature of *Mary* would be co-indexed with the PRO subject of *dance* as desired. There is no way to generate the illegal co-indexing in (63). Thus the locality constraints on PRO control fall out from the mechanics of TAG derivation and feature unification.

8.10 Reported speech

Reported speech is handled in the XTAG grammar by having the reporting clause adjoin into the quote. Thus, the reporting clause is an auxiliary tree, anchored by the reporting verb. See [Doran, 1998] for details of the analysis. There are trees in both the Tnx0Vs1 and Tnx0nx1s2 families to handle reporting clauses which precede, follow and come in the middle of the quote.

Chapter 9

The English Copula, Raising Verbs, and Small Clauses

The English copula, raising verbs, and small clauses are all handled in XTAG by a common analysis based on sentential clauses headed by non-verbal elements. Since there are a number of different analyses in the literature of how these phenomena are related (or not), we will present first the data for all three phenomena, then various analyses from the literature, finishing with the analysis used in the English XTAG grammar.¹

9.1 Usages of the copula, raising verbs, and small clauses

9.1.1 Copula

The verb *be* as used in sentences (65)-(67) is often referred to as the COPULA. It can be followed by a noun, adjective, or prepositional phrase.

(65) Carl is a jerk .

(66) Carl is upset .

(67) Carl is in a foul mood .

Although the copula may look like a main verb at first glance, its syntactic behavior follows the auxiliary verbs rather than main verbs. In particular,

- Copula *be* inverts with the subject.

(68) is Beth writing her dissertation ?
is Beth upset ?
*wrote Beth her dissertation ?

- Copula *be* occurs to the left of the negative marker *not*.

¹This chapter is strongly based on [Heycock, 1991]. Sections 9.1 and 9.2 are greatly condensed from her paper, while the description of the XTAG analysis in section 9.3 is an updated and expanded version.

CHAPTER 9. THE ENGLISH COPULA, RAISING VERBS, AND SMALL CLAUSES

- (69) Beth is not writing her dissertation .
Beth is not upset .
*Beth wrote not her dissertation .

- Copula *be* can contract with the negative marker *not*.

- (70) Beth isn't writing her dissertation .
Beth isn't upset .
*Beth wroten't her dissertation .

- Copula *be* can contract with pronominal subjects.

- (71) She's writing her dissertation .
She's upset .
*She'ote her dissertation .

- Copula *be* occurs to the left of adverbs in the unmarked order.

- (72) Beth is often writing her dissertation .
Beth is often upset .
*Beth wrote often her dissertation .

Unlike all the other auxiliaries, however, copula *be* is not followed by a verbal category (by definition) and therefore must be the rightmost verb. In this respect, it is like a main verb.

The semantic behavior of the copula is also unlike main verbs. In particular, any semantic restrictions or roles placed on the subject come from the complement phrase (NP, AP, PP) rather than from the verb, as illustrated in sentences (73) and (74). Because the complement phrases predicate over the subject, these types of sentences are often called PREDICATIVE sentences.

- (73) The bartender was garrulous .

- (74) ?The cliff was garrulous .

9.1.2 Raising Verbs

Raising verbs are the class of verbs that share with the copula the property that the complement, rather than the verb, places semantic constraints on the subject.

- (75) Carl seems a jerk .
Carl seems upset .
Carl seems in a foul mood .

- (76) Carl appears a jerk .
Carl appears upset .
Carl appears in a foul mood .

The raising verbs are similar to auxiliaries in that they order with other verbs, but they are unique in that they can appear to the left of the infinitive, as seen in the sentences in (77). They cannot, however, invert or contract like other auxiliaries (78), and they appear to the right of adverbs (79).

- (77) Carl seems to be a jerk .
 Carl seems to be upset .
 Carl seems to be in a foul mood .

- (78) *seems Carl to be a jerk ?
 *Carl seemn't to be upset .
 *Carl'ems to be in a foul mood .

- (79) Carl often seems to be upset .
 *Carl seems often to be upset .

9.1.3 Small Clauses

One way of describing small clauses is as predicative sentences without the copula. Since matrix clauses require tense, these clausal structures appear only as embedded sentences. They occur as complements of certain verbs, each of which may allow certain types of small clauses but not others, depending on its lexical idiosyncrasies.

- (80) I consider [Carl a jerk] .
 I consider [Carl upset] .
 ?I consider [Carl in a foul mood] .

- (81) I prefer [Carl in a foul mood] .
 ??I prefer [Carl upset] .

9.1.4 Raising Adjectives

Raising adjectives are the class of adjectives that share with the copula and raising verbs the property that the complement, rather than the verb, places semantic constraints on the subject.

They appear with the copula in a matrix clause, as in (82). However, in other cases, such as that of small clauses (83), they do not have to appear with the copula.

- (82) Carl is likely to be a jerk .
 Carl is likely to be upset .
 Carl is likely to be in a foul mood .
 Carl is likely to perjure himself .
- (83) I consider Carl likely to perjure himself .

9.2 Various Analyses

9.2.1 Main Verb Raising to INFL + Small Clause

In [Pollock, 1989] the copula is generated as the head of a VP, like any main verb such as *sing* or *buy*. Unlike all other main verbs², however, *be* moves out of the VP and into Infl in a

²with the exception of *have* in British English. See footnote 1 in Chapter 20.

tensed sentence. This analysis aims to account for the behavior of *be* as an auxiliary in terms of inversion, negative placement and adverb placement, while retaining a sentential structure in which *be* heads the main VP at D-Structure and can thus be the only verb in the clause.

Pollock claims that the predicative phrase is not an argument of *be*, which instead he assumes to take a small clause complement, consisting of a node dominating an NP and a predicative AP, NP or PP. The subject NP of the small clause then raises to become the subject of the sentence. This accounts for the failure of the copula to impose any selectional restrictions on the subject. Raising verbs such as *seem* and *appear*, presumably, take the same type of small clause complement.

9.2.2 Auxiliary + Null Copula

In [Lapointe, 1980] the copula is treated as an auxiliary verb that takes as its complement a VP headed by a passive verb, a present participle, or a null verb (the true copula). This verb may then take AP, NP or PP complements. The author points out that there are many languages that have been analyzed as having a null copula, but that English has the peculiarity that its null copula requires the co-presence of the auxiliary *be*.

9.2.3 Auxiliary + Predicative Phrase

In GPSG ([Gazdar *et al.*, 1985], [Sag *et al.*, 1985]) the copula is treated as an auxiliary verb that takes an X^2 category with a + value for the head feature [PRD] (predicative). AP, NP, PP and VP can all be [+PRD], but a Feature Co-occurrence Restriction guarantees that a [+PRD] VP will be headed by a verb that is either passive or a present participle.

GPSG follows [Chomsky, 1970] in adopting the binary valued features [V] and [N] for decomposing the verb, noun, adjective and preposition categories. In that analysis, verbs are [+V, −N], nouns are [−V, +N], adjectives [+V, +N] and prepositions [−V, −N]. NP and AP predicative complements generally pattern together; a fact that can be stated economically using this category decomposition. In neither [Sag *et al.*, 1985] nor [Chomsky, 1970] is there any discussion of how to handle the complete range of complements to a verb like *seem*, which takes AP, NP and PP complements, as well as infinitives. The solution would appear to be to associate the verb with two sets of rules for small clauses, leaving aside the use of the verb with an expletive subject and sentential complement.

9.2.4 Auxiliary + Small Clause

In [Moro, 1990] the copula is treated as a special functional category - a lexicalization of tense, which is considered to head its own projection. It takes as a complement the projection of another functional category, Agr (agreement). This projection corresponds roughly to a small clause, and is considered to be the domain within which predication takes place. An NP must then raise out of this projection to become the subject of the sentence: it may be the subject of the AgrP, or, if the predicate of the AgrP is an NP, this may raise instead. In addition to occurring as the complement of *be*, AgrP is selected by certain verbs such as *consider*. It follows from this analysis that when the complement to *consider* is a simple AgrP, it will always consist of a subject followed by a predicate, whereas if the complement contains the verb *be*,

the predicate of the AgrP may raise to the left of *be*, leaving the subject of the AgrP to the right.

- (84) John_i is [_{AgrP} *t_i* the culprit] .
 (85) The culprit_i is [_{AgrP} John *t_i*] .
 (86) I consider [_{AgrP} John the culprit] .
 (87) I consider [John_i to be [_{AgrP} *t_i* the culprit]] .
 (88) I consider [the culprit_i to be [_{AgrP} John *t_i*]] .

Moro does not discuss a number of aspects of his analysis, including the nature of Agr and the implied existence of sentences without VP's.

9.3 XTAG analysis

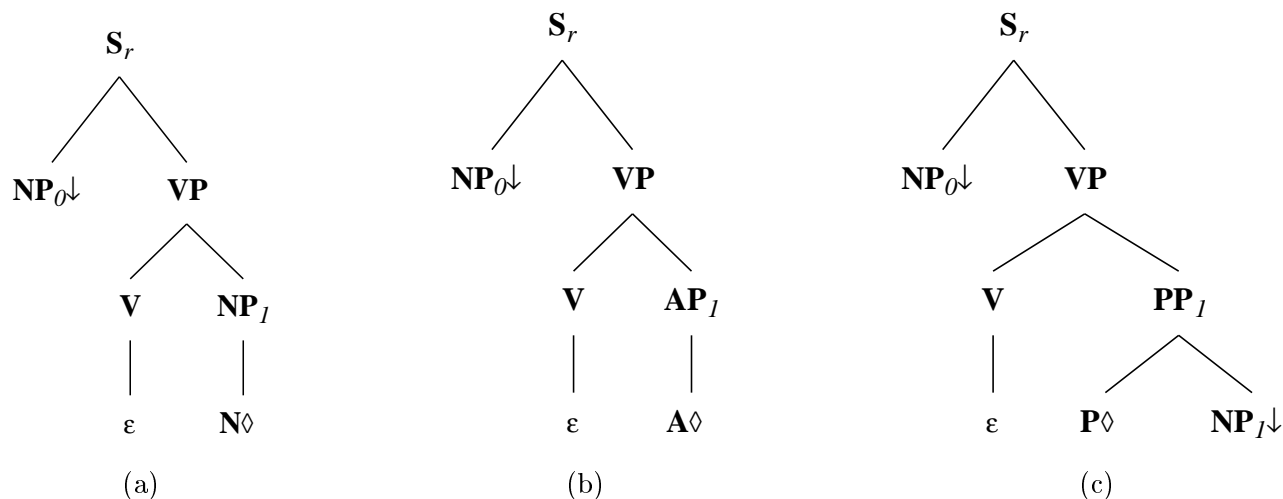


Figure 9.1: Predicative trees: $\alpha nx0N1$ (a), $\alpha nx0Ax1$ (b) and $\alpha nx0Pnx1$ (c)

The XTAG grammar provides a uniform analysis for the copula, raising verbs and small clauses by treating the maximal projections of lexical items that can be predicated as predicative clauses, rather than simply noun, adjective and prepositional phrases. The copula adjoins in for matrix clauses, as do the raising verbs. Certain other verbs (such as *consider*) can take the predicative clause as a complement, without the adjunction of the copula, to form the embedded small clause.

The structure of a predicative clause, then, is roughly as seen in (89)-(91) for NP's, AP's and PP's. The XTAG trees corresponding to these structures³ are shown in Figures 9.1(a), 9.1(b), and 9.1(c), respectively.

³There are actually two other predicative trees in the XTAG grammar. Another predicative noun phrase tree is needed for noun phrases without determiners, as in the sentence *They are firemen*, and another prepositional phrase tree is needed for exhaustive prepositional phrases, such as *The workers are below*.

CHAPTER 9. THE ENGLISH COPULA, RAISING VERBS, AND SMALL CLAUSES

(89) [_S NP [_{VP} N ...]]

(90) [_S NP [_{VP} A ...]]

(91) [_S NP [_{VP} P ...]]

The copula *be* and raising verbs all get the basic auxiliary tree as explained in the section on auxiliary verbs (section 20.1). Unlike the raising verbs, the copula also selects the inverted auxiliary tree set. Figure 9.2 shows the basic auxiliary tree anchored by the copula *be*. The **<mode>** feature is used to distinguish the predicative constructions so that only the copula and raising verbs adjoin onto the predicative trees.

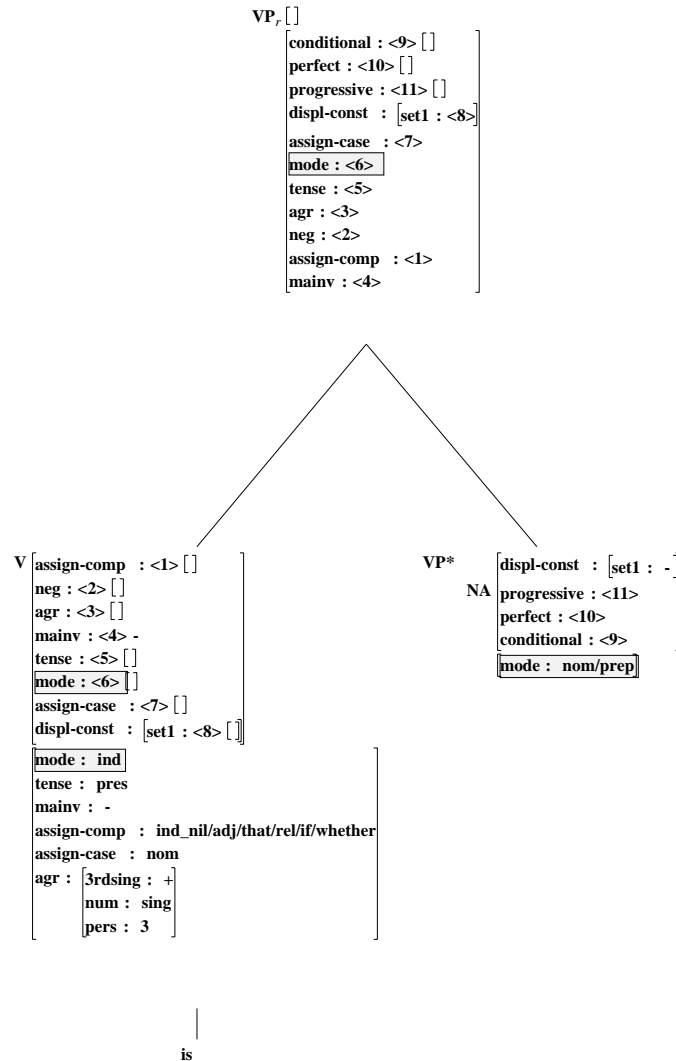


Figure 9.2: Copula auxiliary tree: βV_{vx}

There are two possible values of **<mode>** that correspond to the predicative trees, **nom** and **prep**. They correspond to a modified version of the four-valued [N,V] feature described

in section 9.2.3. The **nom** value corresponds to [N+], selecting the NP and AP predicative clauses. As mentioned earlier, they often pattern together with respect to constructions using predicative clauses. The remaining prepositional phrase predicative clauses, then, correspond to the **prep** mode.

Figure 9.3 shows the predicative adjective tree from Figure 9.1(b) now anchored by *upset* and with the features visible. As mentioned, $\langle \mathbf{mode} \rangle = \mathbf{nom}$ on the VP node prevents auxiliaries other than the copula or raising verbs from adjoining into this tree. In addition, it prevents the predicative tree from occurring as a matrix clause. Since all matrix clauses in XTAG must be mode indicative (**ind**) or imperative (**imp**), a tree with $\langle \mathbf{mode} \rangle = \mathbf{nom}$ or $\langle \mathbf{mode} \rangle = \mathbf{prep}$ must have an auxiliary verb (the copula or a raising verb) adjoin in to make it $\langle \mathbf{mode} \rangle = \mathbf{ind}$.

The distribution of small clauses as embedded complements to some verbs is also managed through the mode feature. Verbs such as *consider* and *prefer* select trees that take a sentential complement, and then restrict that complement to be $\langle \mathbf{mode} \rangle = \mathbf{nom}$ and/or $\langle \mathbf{mode} \rangle = \mathbf{prep}$, depending on the lexical idiosyncrasies of that particular verb. Many verbs that don't take small clause complements do take sentential complements that are $\langle \mathbf{mode} \rangle = \mathbf{ind}$, which includes small clauses with the copula already adjoined. Hence, as seen in sentence sets (92)-(94), *consider* takes only small clause complements, *prefer* takes both **prep** (but not **nom**) small clauses and indicative clauses, while *feel* takes only indicative clauses.

- (92) She considers Carl a jerk .
 ?She considers Carl in a foul mood .
 *She considers that Carl is a jerk .
- (93) *She prefers Carl a jerk .
 She prefers Carl in a foul mood .
 She prefers that Carl is a jerk .
- (94) *She feels Carl a jerk .
 *She feels Carl in a foul mood .
 She feels that Carl is a jerk .

Figure 9.4 shows the tree anchored by *consider* that takes the predicative small clauses.

Raising verbs such as *seems* work essentially the same as the auxiliaries, in that they also select the basic auxiliary tree, as in Figure 9.2. The only difference is that the value of $\langle \mathbf{mode} \rangle$ on the VP foot node might be different, depending on what types of complements the raising verb takes. Also, two of the raising verbs take an additional tree, βV_{pxvx} , shown in Figure 9.5, which allows for an experiencer argument, as in *John seems to me to be happy*.

Raising adjectives, such as *likely*, take the tree shown in Figure 9.6. This tree combines aspects of the auxiliary tree βV_{vx} and the adjectival predicative tree shown in Figure 9.1(b). As with βV_{vx} , it adjoins in as a VP auxiliary tree. However, since it is anchored by an adjective, not a verb, it is similar to the adjectival predicative tree in that it has an ϵ at the V node, and a feature value of $\langle \mathbf{mode} \rangle = \mathbf{nom}$ which is passed up to the VP root indicates that it is an adjectival predication. This serves the same purpose as in the case of the tree in Figure 9.3, and forces another auxiliary verb, such as the copula, to adjoin in to make it $\langle \mathbf{mode} \rangle = \mathbf{ind}$.

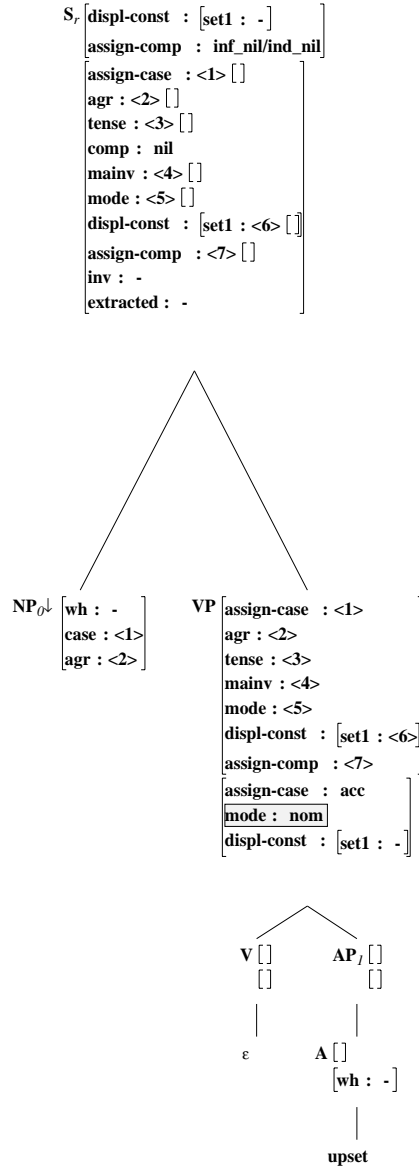


Figure 9.3: Predicative AP tree with features: $\alpha n x 0 A x 1$

9.4 Non-predicative *BE*

The examples with the copula that we have given seem to indicate that *be* is always followed by a predicative phrase of some sort. This is not the case, however, as seen in sentences such as (95)-(100). The noun phrases in these sentences are not predicative. They do not take raising verbs, and they do not occur in embedded small clause constructions.

(95) my teacher is Mrs. Wayman .

(96) Doug is the man with the glasses .

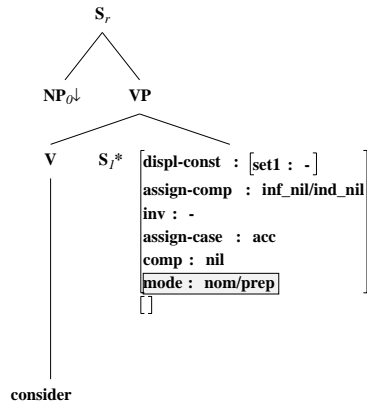


Figure 9.4: *Consider* tree for embedded small clauses

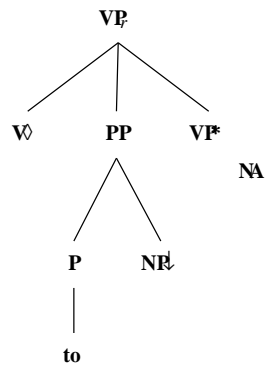


Figure 9.5: Raising verb with experiencer tree: βV_{pxvx}

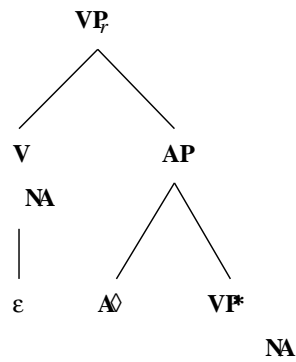


Figure 9.6: Raising adjective tree: βV_{vx-adj}

(97) *My teacher seems Mrs. Wayman .

(98) *Doug appears the man with the glasses .

(99) *I consider [my teacher Mrs. Wayman] .

(100) *I prefer [Doug the man with the glasses] .

In addition, the subject and complement can exchange positions in these type of examples but not in sentences with predicative *be*. Sentence (101) has the same interpretation as sentence (96) and differs only in the positions of the subject and complement NP's. Similar sentences, with a predicative *be*, are shown in (102) and (103). In this case, the sentence with the exchanged NP's (103) is ungrammatical.

(101) The man with the glasses is Doug .

(102) Doug is a programmer .

(103) *A programmer is Doug .

The non-predicative *be* in (95) and (96), also called EQUATIVE BE, patterns differently, both syntactically and semantically, from the predicative usage of *be*. Since these sentences are clearly not predicative, it is not desirable to have a tree structure that is anchored by the NP, AP, or PP, as we have in the predicative sentences. In addition to the conceptual problem, we would also need a mechanism to block raising verbs from adjoining into these sentences (while allowing them for true predicative phrases), and prevent these types of sentence from being embedded (again, while allowing them for true predicative phrases).

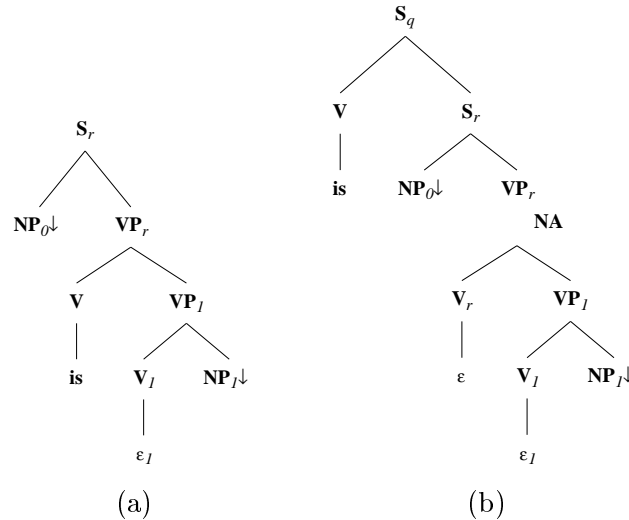


Figure 9.7: Equative *BE* trees: $\alpha_{nx0BEnx1}$ (a) and $\alpha_{Invnx0BEnx1}$ (b)

Although non-predicative *be* is not a raising verb, it does exhibit the auxiliary verb behavior set out in section 9.1.1. It inverts, contracts, and so forth, as seen in sentences (104) and (105), and therefore can not be associated with any existing tree family for main verbs. It requires a separate tree family that includes the tree for inversion. Figures 9.7(a) and 9.7(b) show the declarative and inverted trees, respectively, for equative *be*.

(104) is my teacher Mrs. Wayman ?

(105) Doug isn't the man with the glasses .

Chapter 10

Ditransitive constructions and dative shift

Verbs such as *give* and *put* that require two objects, as shown in examples (106)-(109), are termed ditransitive.

(106) Christy gave a cannoli to Beth Ann .

(107) *Christy gave Beth Ann .

(108) Christy put a cannoli in the refrigerator .

(109) *Christy put a cannoli .

The indirect objects *Beth Ann* and *refrigerator* appear in these examples in the form of PP's. Within the set of ditransitive verbs there is a subset that also allow two NP's as in (110). As can be seen from (110) and (111) this two NP, or double-object, construction is grammatical for *give* but not for *put*.

(110) Christy gave Beth Ann a cannoli .

(111) *Christy put the refrigerator the cannoli .

The alternation between (106) and (110) is known as dative shift.¹ In order to account for verbs with dative shift the English XTAG grammar includes structures for both variants in the tree family $T_{nx0V_{nx1}P_{nx2}}$. The declarative trees for the shifted and non-shifted alternations are shown in Figure 10.1.

The indexing of nodes in these two trees represents the fact that the semantic role of the indirect object (NP_2) in Figure 10.1(a) is the same as that of the direct object (NP_2) in Figure 10.1(b) (and vice versa). This use of indexing is consistent with our treatment of other constructions such as passive and ergative.

Verbs that do not show this alternation and have only the NP PP structure (e.g. *put*) select the tree family $T_{nx0V_{nx1}P_{nx2}}$. Unlike the $T_{nx0V_{nx1}P_{nx2}}$ family, the $T_{nx0V_{nx1}P_{nx2}}$ tree family does not contain trees for the NP NP structure. Other verbs such as *ask* allow only the NP NP structure as shown in (112) and (113).

¹In languages similar to English that have overt case marking indirect objects would be marked with dative case. It has also been suggested that for English the preposition *to* serves as a dative case marker.

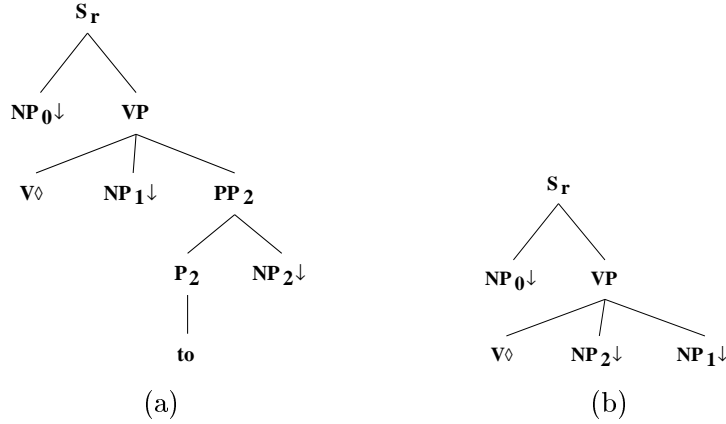


Figure 10.1: Dative shift trees: $\alpha nx0Vnx1Pnx2$ (a) and $\alpha nx0Vnx2nx1$ (b)

(112) Beth Ann asked Sridhar a question .

(113) *Beth Ann asked a question to Sridhar .

Verbs that only allow the NP NP structure select the tree family $Tnx0Vnx1nx2$. This tree family does not have the trees for the NP PP structure.

Notice that in Figure 10.1(a) the preposition *to* is built into the tree. There are other apparent cases of dative shift with *for*, such as in (114) and (115), that we have taken to be structurally distinct from the cases with *to*.

(114) Beth Ann baked Dusty a biscuit .

(115) Beth Ann baked a biscuit for Dusty .

[McCawley, 1988] notes:

A “*for-dative*” expression in underlying structure is external to the V with which it is combined, in view of the fact that the latter behaves as a unit with regard to all relevant syntactic phenomena.

In other words, the *for* PP’s that appear to undergo dative shift are actually adjuncts, not complements. Examples (116) and (117) demonstrate that PP’s with *for* are optional while ditransitive *to* PP’s are not.

(116) Beth Ann made dinner .

(117) *Beth Ann gave dinner .

Consequently, in the XTAG grammar the apparent dative shift with *for* PP’s is treated as $Tnx0Vnx1nx2$ for the NP NP structure, and as a transitive plus an adjoined adjunct PP for the NP PP structure. To account for the ditransitive *to* PP’s, the preposition *to* is built into the

tree family Tnx0Vnx1tonx2. This accounts for the fact that *to* is the only preposition allowed in dative shift constructions.

[McCawley, 1988] also notes that the *to* and *for* cases differ with respect to passivization; the indirect objects with *to* may be the subjects of corresponding passives while the alleged indirect objects with *for* cannot, as in sentences (118)-(121). Note that the passivisation examples are for NP NP structures of verbs that take *to* or *for* PP's.

(118) Beth Ann gave Clove dinner .

(119) Clove was given dinner (by Beth Ann) .

(120) Beth Ann made Clove dinner .

(121) ?Clove was made dinner (by Beth Ann) .

However, we believe that this to be incorrect, and that the indirect objects in the *for* case are allowed to be the subjects of passives, as in sentences (122)-(123). The apparent strangeness of sentence (121) is caused by interference from other interpretations of *Clove was made dinner* .

(122) Dania baked Doug a cake .

(123) Doug was baked a cake by Dania .

Chapter 11

It-clefts

There are several varieties of it-clefts in English. All the it-clefts have four major components:

- **the dummy subject:** *it*,
- **the main verb:** *be*,
- **the clefted element:** A constituent (XP) compatible with any gap in the clause,
- **the clause:** A clause (e.g. S) with or without a gap.

Examples of it-clefts are shown in (124)-(127).

(124) it was [_{XP} here _{XP}] [_S that the ENIAC was created . _S]

(125) it was [_{XP} at MIT _{XP}] [_S that colorless green ideas slept furiously . _S]

(126) it is [_{XP} happily _{XP}] [_S that Seth quit Reality . _S]

(127) it was [_{XP} there _{XP}] [_S that she would have to enact her renunciation . _S]

The clefted element can be of a number of categories, for example NP, PP or adverb. The clause can also be of several types. The English XTAG grammar currently has a separate analysis for only a subset of the ‘specificational’ it-clefts¹, in particular the ones without gaps in the clause (e.g. (126) and (127)). It-clefts that have gaps in the clause, such as (124) and (125) are currently handled as relative clauses. Although arguments have been made against treating the clefted element and the clause as a constituent ([Delahunty, 1984]), the relative clause approach does capture the restriction that the clefted element must fill the gap in the clause, and does not require any additional trees.

In the ‘specificational’ it-cleft without gaps in the clause, the clefted element has the role of an adjunct with respect to the clause. For these cases the English XTAG grammar requires additional trees. These it-cleft trees are in separate tree families because, although some researchers (e.g. [Akmajian, 1970]) derived it-clefts through movement from other sentence types, most current researchers (e.g. [Delahunty, 1984], [Knowles, 1986], [Gazdar *et al.*, 1985], [Delin,

¹See e.g. [Ball, 1991], [Delin, 1989] and [Delahunty, 1984] for more detailed discussion of types of it-clefts.

1989] and [Sornicola, 1988]) favor base-generation of the various cleft sentences. Placing the it-cleft trees in their own tree families is consistent with the current preference for base generation, since in the XTAG English grammar, structures that would be related by transformation in a movement-based account will appear in the same tree family. Like the base-generated approaches, the placement of it-clefts in separate tree families makes the claim that there is no derivational relation between it-clefts and other sentence types.

The three it-cleft tree families are virtually identical except for the category label of the clefted element. Figure 11.1 shows the declarative tree and an inverted tree for the PP It-cleft tree family.

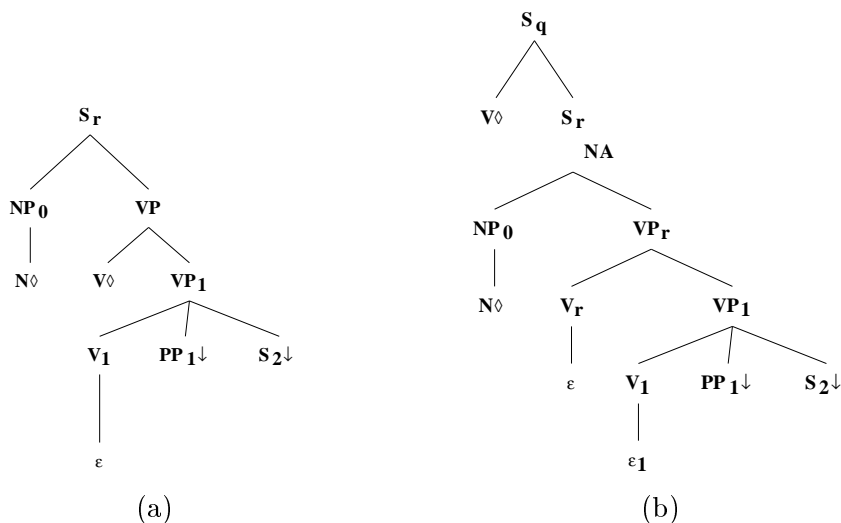


Figure 11.1: It-cleft with PP clefted element: $\alpha\text{ItVpnx1s2}$ (a) and $\alpha\text{InvItVpnx1s2}$ (b)

The extra layer of tree structure in the VP represents that, while *be* is a main verb rather than an auxiliary in these cases, it retains some auxiliary properties. The VP structure for the equative/it-cleft-*be* is identical to that obtained after adjunction of predicative-*be* into small-clauses.² The inverted tree in Figure 11.1(b) is necessary because of *be*'s auxiliary-like behavior. Although *be* is the main verb in it-clefts, it inverts like an auxiliary. Main verb inversion cannot be accomplished by adjunction as is done with auxiliaries and therefore must be built into the tree family. The tree in Figure 11.1(b) is used for yes/no questions such as (128).

(128) was it in the forest that the wolf talked to the little girl ?

²For additional discussion of equative or predicative-*be* see Chapter 9.

Part III

Sentence Types

Chapter 12

Passives

In passive constructions such as (129), the subject NP is interpreted as having the same role as the direct object NP in the corresponding active declarative (130).

(129) **An airline buy-out bill** was approved by the House. (WSJ)

(130) The House approved **an airline buy-out bill**.

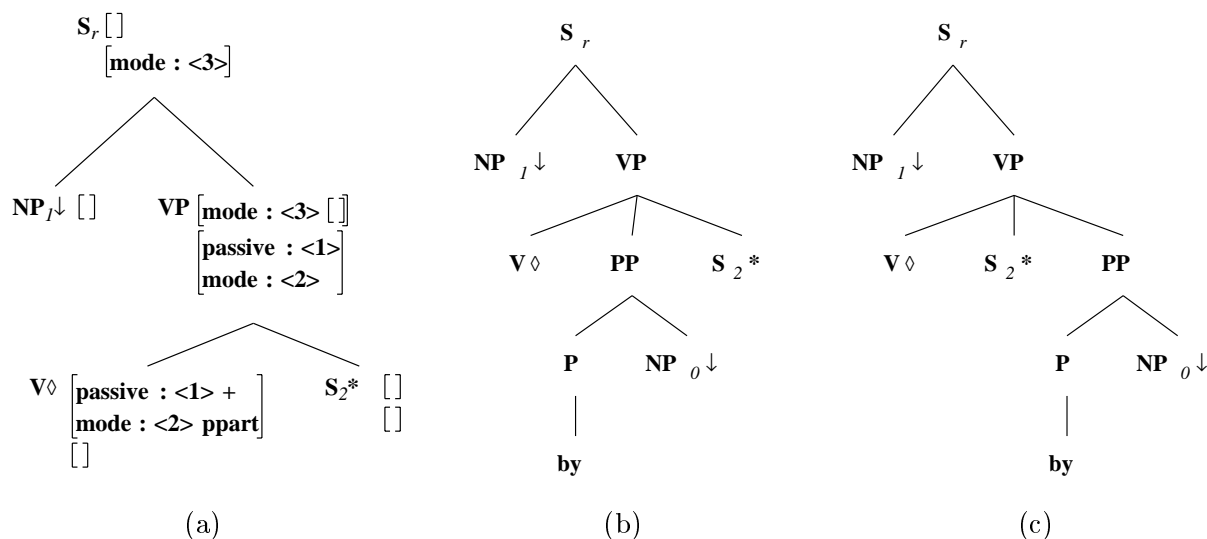


Figure 12.1: Passive trees in the Sentential Complement with NP tree family: β_{nx1Vs2} (a), $\beta_{nx1Vbynxs2}$ (b) and $\beta_{nx1Vs2bynxs0}$ (c)

In a movement analysis, the direct object is said to have moved to the subject position. The original declarative subject is either absent in the passive or is in a *by* headed PP (*by* phrase). In the English XTAG grammar, passive constructions are handled by having separate trees within the appropriate tree families. Passive trees are found in most tree families that have a direct object in the declarative tree (the light verb tree families, for instance, do not contain passive trees). Passive trees occur in pairs - one tree with the *by* phrase, and another without

it. Variations in the location of the *by* phrase are possible if a subcategorization includes other arguments such as a PP or an indirect object. Additional trees are required for these variations. For example, the Sentential Complement with NP tree family has three passive trees, shown in Figure 12.1: one without the *by*-phrase (Figure 12.1(a)), one with the *by* phrase before the sentential complement (Figure 12.1(b)), and one with the *by* phrase after the sentential complement (Figure 12.1(c)).

Figure 12.1(a) also shows the feature restrictions imposed on the anchor¹. Only verbs with **<mode>=ppart** (i.e. verbs with passive morphology) can anchor this tree. The **<mode>** feature is also responsible for requiring that passive *be* adjoin into the tree to create a matrix sentence. Since a requirement is imposed that all matrix sentences must have **<mode>=ind/imp**, an auxiliary verb that selects **<mode>=ppart** and **<passive>=+** (such as *was*) must adjoin (see Chapter 20 for more information on the auxiliary verb system).

¹A reduced set of features are shown for readability.

Chapter 13

Extraction

The discussion in this chapter covers constructions that are analyzed as having wh-movement in GB, in particular, wh-questions and topicalization. Relative clauses, which could also be considered extractions, are discussed in Chapter 14.

Extraction involves a constituent appearing in a linear position to the left of the clause with which it is interpreted. One clause argument position is empty. For example, the position filled by *frisbee* in the declarative in sentence (131) is empty in sentence (132). The wh-item *what* in sentence (132) is of the same syntactic category as *frisbee* in sentence (131) and fills the same role with respect to the subcategorization.

(131) Clove caught a frisbee.

(132) What_{*i*} did Clove catch ϵ_i ?

The English XTAG grammar represents the connection between the extracted element and the empty position with co-indexing (as does GB). The <**trace**> feature is used to implement the co-indexing. In extraction trees in XTAG, the ‘empty’ position is filled with an ϵ . The extracted item always appears in these trees as a sister to the S_r tree, with both dominated by a S_q root node. The S_r subtrees in extraction trees have the same structure as the declarative tree in the same tree family. The additional structure in extraction trees of the S_q and NP nodes roughly corresponds to the CP and Spec of CP positions in GB.

All sentential trees with extracted components (this does not include relative clause trees) are marked <**extracted**>=+ at the top S node, while sentential trees with no extracted components are marked <**extracted**>=-. Items that take embedded sentences, such as nouns, verbs and some prepositions can place restrictions on whether the embedded sentence is allowed to be extracted or not. For instance, sentential subjects and sentential complements of nouns and prepositions are not allowed to be extracted, while certain verbs may allow extracted sentential complements and others may not (e.g. sentences (133)-(136)).

(133) The jury wondered [who killed Nicole].

(134) The jury wondered [who Simpson killed].

(135) The jury thought [Simpson killed Nicole].

(136) *The jury thought [who did Simpson kill]?

The **<extracted>** feature is also used to block embedded topicalization in infinitival complement clauses as exemplified in (137).

(137) * John wants [Bill_i [PRO to see t_i]]

Verbs such as *want* that take non-*wh* infinitival complements specify that the **<extracted>** feature of their complement clause (i.e. of the foot S node) is $-$. Clauses that involve topicalization have $+$ as the value of their **<extracted>** feature (i.e. of the root S node). Sentences like (137) are thus ruled out.

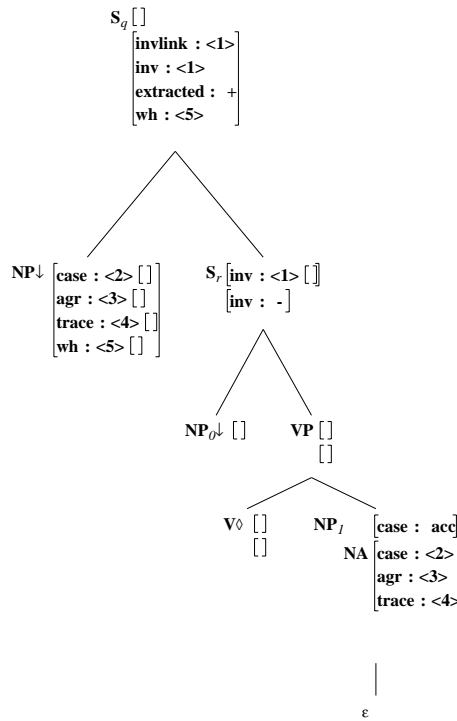


Figure 13.1: Transitive tree with object extraction: $\alpha W1nx0Vnx1$

The tree that is used to derive the embedded sentence in (135) in the English XTAG grammar is shown in Figure 13.1¹. The important features of extracted trees are:

- The subtree that has S_r as its root is identical to the declarative tree or a non-extracted passive tree, except for having one NP position in the VP filled by ϵ .
- The root S node is S_q , which dominates NP and S_r .
- The **<trace>** feature of the ϵ filled NP is co-indexed with the **<trace>** feature of the NP daughter of S_q .

¹Features not pertaining to this discussion have been taken out to improve readability.

- The **<case>** and **<agr>** features are passed from the empty NP to the extracted NP. This is particularly important for extractions from subject NP's, since **<case>** can continue to be assigned from the verb to the subject NP position, and from there be passed to the extracted NP.
- The **<inv>** feature of S_r is co-indexed to the **<wh>** feature of NP through the use of the **<invlink>** feature in order to force subject-auxiliary inversion where needed (see section 13.1 for more discussion of the **<inv>**/**<wh>** co-indexing and the use of these trees for topicalization).

13.1 Topicalization and the value of the **<inv>** feature

Our analysis of topicalization uses the same trees as wh-extraction. For any NP complement position a single tree is used for both wh-questions and for topicalization from that position. Wh-questions have subject-auxiliary inversion and topicalizations do not. This difference between the constructions is captured by equating the values of the S_r 's **<inv>** feature and the extracted NP's **<wh>** feature. This means that if the extracted item is a wh-expression, as in wh-questions, the value of **<inv>** will be + and an inverted auxiliary will be forced to adjoin. If the extracted item is a non-wh, **<inv>** will be – and no auxiliary adjunction will occur. An additional complication is that inversion only occurs in matrix clauses, so the values of **<inv>** and **<wh>** should only be equated in matrix clauses and not in embedded clauses. In the English XTAG grammar, appropriate equating of the **<inv>** and **<wh>** features is accomplished using the **<invlink>** feature and a restriction imposed on the root S of a derivation. In particular, in extraction trees that are used for both wh-questions and topicalizations, the value of the **<inv>** feature for the top of the S_r node is co-indexed to the value of the **<inv>** feature on the bottom of the S_q node. On the bottom of the S_q node the **<inv>** feature is co-indexed to the **<invlink>** feature. The **<wh>** feature of the extracted NP node is co-indexed to the value of the **<wh>** feature on the bottom of S_q . The linking between the value of the S_q **<wh>** and the **<invlink>** features is imposed by a condition on the final root node of a derivation (i.e. the top S node of a matrix clause) requires that **<invlink>=<wh>**. For example, the tree in Figure 13.1 is used to derive both (138) and (139).

(138) John, I like.

(139) Who do you like?

For the question in (139), the extracted item *who* has the feature value **<wh>=+**, so the value of the **<inv>** feature on VP is also + and an auxiliary, in this case *do*, is forced to adjoin. For the topicalization (138) the values for *John's* **<wh>** feature and for S_q 's **<inv>** feature are both – and no auxiliary adjoins.

13.2 Extracted subjects

The extracted subject trees provide for sentences like (140)-(142), depending on the tree family with which it is associated.

(140) Who left?

(141) Who wrote the paper?

(142) Who was happy?

Wh-questions on subjects differ from other argument extractions in not having subject-auxiliary inversion. This means that in subject wh-questions the linear order of the constituents is the same as in declaratives so it is difficult to tell whether the subject has moved out of position or not (see [Heycock and Kroch, 1993] for arguments for and against moved subject).

The English XTAG treatment of subject extractions assumes the following:

- Syntactic subject topicalizations don't exist; and
- Subjects in wh-questions are extracted rather than in situ.

The assumption that there is no syntactic subject topicalization is reasonable in English since there is no convincing syntactic evidence and since the interpretability of subjects as topics seems to be mainly affected by discourse and intonational factors rather than syntactic structure. As for the assumption that wh-question subjects are extracted, these questions seem to have more similarities to other extractions than to the two cases in English that have been considered in situ wh: multiple wh questions and echo questions. In multiple wh questions such as sentence (143), one of the wh-items is blocked from moving sentence initially because the first wh-item already occupies the location to which it would move.

(143) Who ate what?

This type of 'blocking' account is not applicable to subject wh-questions because there is no obvious candidate to do the blocking. Similarity between subject wh-questions and echo questions is also lacking. At least one account of echo questions ([Hockey, 1994]) argues that echo questions are not ordinary wh-questions at all, but rather focus constructions in which the wh-item is the focus. Clearly, this is not applicable to subject wh-questions. So it seems that treating subject wh-questions similarly to other wh-extractions is more justified than an in situ treatment.

Given these assumptions, there must be separate trees in each tree family for subject extractions. The declarative tree cannot be used even though the linear order is the same because the structure is different. Since topicalizations are not allowed, the $\langle \mathbf{wh} \rangle$ feature for the extracted NP node is set in these trees to $+$. The lack of subject-auxiliary inversion is handled by the absence of the $\langle \mathbf{invlink} \rangle$ feature. Without the presence of this feature, the $\langle \mathbf{wh} \rangle$ and $\langle \mathbf{inv} \rangle$ are never linked, so inversion can not occur. Like other wh-extractions, the S_q node is marked $\langle \mathbf{extracted} \rangle = +$ to constrain the occurrence of these trees in embedded sentences. The tree in Figure 13.2 is an example of a subject wh-question tree.

13.3 Wh-moved NP complement

Wh-questions can be formed on every NP object or indirect object that appears in the declarative tree or in the passive trees, as seen in sentences (144)-(149). A tree family will contain

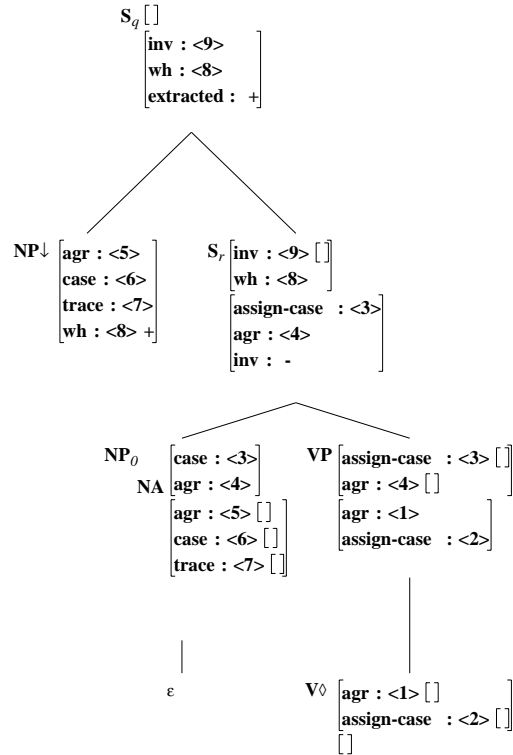


Figure 13.2: Intransitive tree with subject extraction: $\alpha W0nx0V$

one tree for each of these possible NP complement positions. Figure 13.3 shows the two extraction trees from the ditransitive tree family for the extraction of the direct (Figure 13.3(a)) and indirect object (Figure 13.3(b)).

- (144) Dania asked Beth a question.
(145) Who_i did Dania ask ϵ_i a question?
(146) What_i did Dania ask Beth ϵ_i ?
(147) Beth was asked a question by Dania.
(148) Who_i was Beth asked a question by ϵ_i ??
(149) What_i was Beth asked ϵ_i ? by Dania?

13.4 Wh-moved object of a P

Wh-questions can be formed on the NP object of a complement PP as in sentence (150).

- (150) [Which dog]_i did Beth Ann give a bone to ϵ_i ?

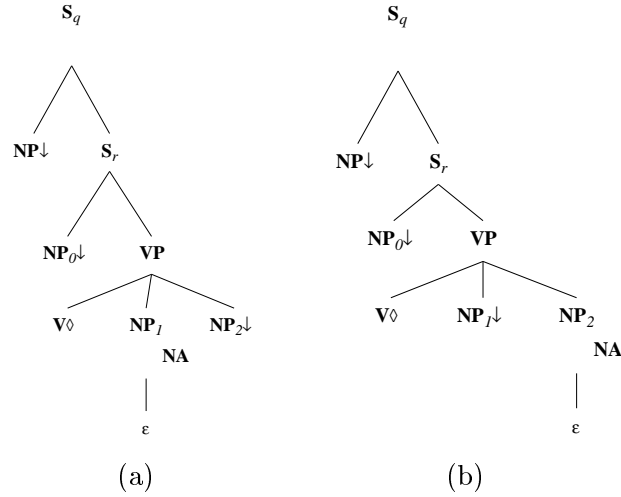


Figure 13.3: Ditransitive trees with direct object: $\alpha W1nx0Vnx1nx2$ (a) and indirect object extraction: $\alpha W2nx0Vnx1nx2$ (b)

The *by* phrases of passives behave like complements and can undergo the same type of extraction, as in (151).

(151) [Which dog]_{*i*} was the frisbee caught by ϵ_i ?

Tree structures for this type of sentence are very similar to those for the *wh*-extraction of NP complements discussed in section 13.3 and have the identical important features related to tree structure and trace and inversion features. The tree in Figure 13.4 is an example of this type of tree. Topicalization of NP objects of prepositions is handled the same way as topicalization of complement NP's.

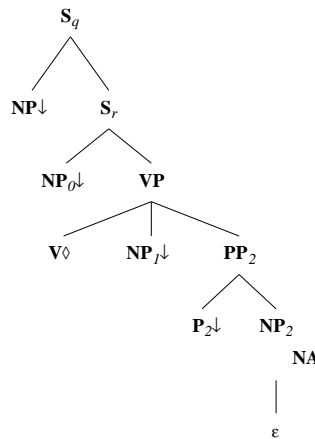


Figure 13.4: Ditransitive with PP tree with the object of the PP extracted: $\alpha W2nx0Vnx1pnx2$

13.5 Wh-moved PP

Like NP complements, PP complements can be extracted to form wh-questions, as in sentence (152).

(152) [To which dog]_i did Beth Ann throw the frisbee ϵ_i ?

As can be seen in the tree in Figure 13.5, extraction of PP complements is very similar to extraction of NP complements from the same positions.

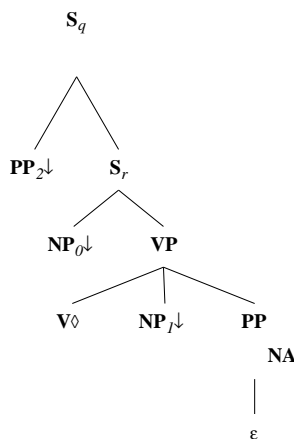


Figure 13.5: Ditransitive with PP with PP extraction tree: $\alpha pW2nx0Vnx1pnx2$

The PP extraction trees differ from NP extraction trees in having a PP rather than an NP left daughter node under S_q and in having the ϵ fill a PP rather than an NP position in the VP. In other respects these PP extraction structures behave like the NP extractions, including being used for topicalization.

13.6 Wh-moved S complement

Except for the node label on the extracted position, the trees for wh-questions on S complements look exactly like the trees for wh-questions on NP's in the same positions. This is because there is no separate wh-lexical item for clauses in English, so the item *what* is ambiguous between representing a clause or an NP. To illustrate this ambiguity notice that the question in (153) could be answered by either a clause as in (154) or an NP as in (155). The extracted NP in these trees is constrained to be $\langle \mathbf{wh} \rangle = +$, since sentential complements can not be topicalized.

(153) What does Clove want?

(154) for Beth Ann to play frisbee with her

(155) a biscuit

13.7 Wh-moved Adjective complement

In subcategorizations that select an adjective complement, that complement can be questioned in a wh-question, as in sentence (156).

(156) How_{*i*} did he feel ϵ_i ?

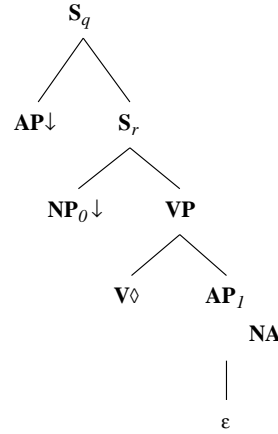


Figure 13.6: Predicative Adjective tree with extracted adjective: $\alpha\text{WA1nx0Vax1}$

The tree families with adjective complements include trees for such adjective extractions that are very similar to the wh-extraction trees for other categories of complements. The adjective position in the VP is filled by an ϵ and the trace feature of the adjective complement and of the adjective daughter of S_q are co-indexed. The extracted adjective is required to be $\langle \mathbf{wh} \rangle = +^2$, so no topicalizations are allowed. An example of this type of tree is shown in Figure 13.6.

²How is the only $\langle \mathbf{wh} \rangle = +$ adjective currently in the XTAG English grammar.

Chapter 14

Relative Clauses

Relative clauses are NP modifiers, which involve extraction of an argument or an adjunct. The NP head (the portion of the NP being modified by the relative clause) is not directly related to the extracted element. For example in (157), *the person* is the head NP and is modified by the relative clause *whose mother ϵ likes Chris*. *The person* is not interpreted as the subject of the relative clause which is missing an overt subject. In other cases, such as (158), the relationship between the head NP *export exhibitions* may seem to be more direct but even there we assume that there are two independent relationships: one between the entire relative clause and the NP it modifies, and another between the extracted element and its trace. The extracted element may be an overt *wh*-phrase as in (157) or a covert element as in (158).

(157) the person whose mother likes Chris

(158) export exhibitions that included high-tech items

Relative clauses are represented in the English XTAG grammar by auxiliary trees that adjoin to NP's. These trees are anchored by the verb in the clause and appear in the appropriate tree families for the various verb subcategorizations. Within a tree family there will be groups of relative clause trees based on the declarative tree and each passive tree. Within each of these groups, there is a separate relative clause tree corresponding to each possible argument that can be extracted from the clause. There is no relationship between the extracted position and the head NP. The relationship between the relative clause and the head NP is treated as a semantic relationship which will be provided by any reasonable compositional theory. The relationship between the extracted element (which can be covert) is captured by co-indexing the **<trace>** features of the extracted NP and the NP_w node in the relative clause tree. If for example, it is NP₀ that is extracted, we have the following feature equations:

$$\text{NP}_{w.t}:\langle \text{trace} \rangle = \text{NP}_0.t:\langle \text{trace} \rangle$$

$$\text{NP}_{w.t}:\langle \text{case} \rangle = \text{NP}_0.t:\langle \text{case} \rangle$$

$$\text{NP}_{w.t}:\langle \text{agr} \rangle = \text{NP}_0.t:\langle \text{agr} \rangle^1$$

Representative examples from the transitive tree family are shown with a relevant subset of their features in Figures 14.1(a) and 14.1(b). Figure 14.1(a) involves a relative clause with a covert extracted element, while figure 14.1(b) involves a relative clause with an overt *wh*-phrase.²

¹No adjunct traces are represented in the XTAG analysis of adjunct extraction. Relative clauses on adjuncts do not have traces and consequently feature equations of the kind shown here are not present.

²The convention followed in naming relative clause trees is outlined in Appendix D.

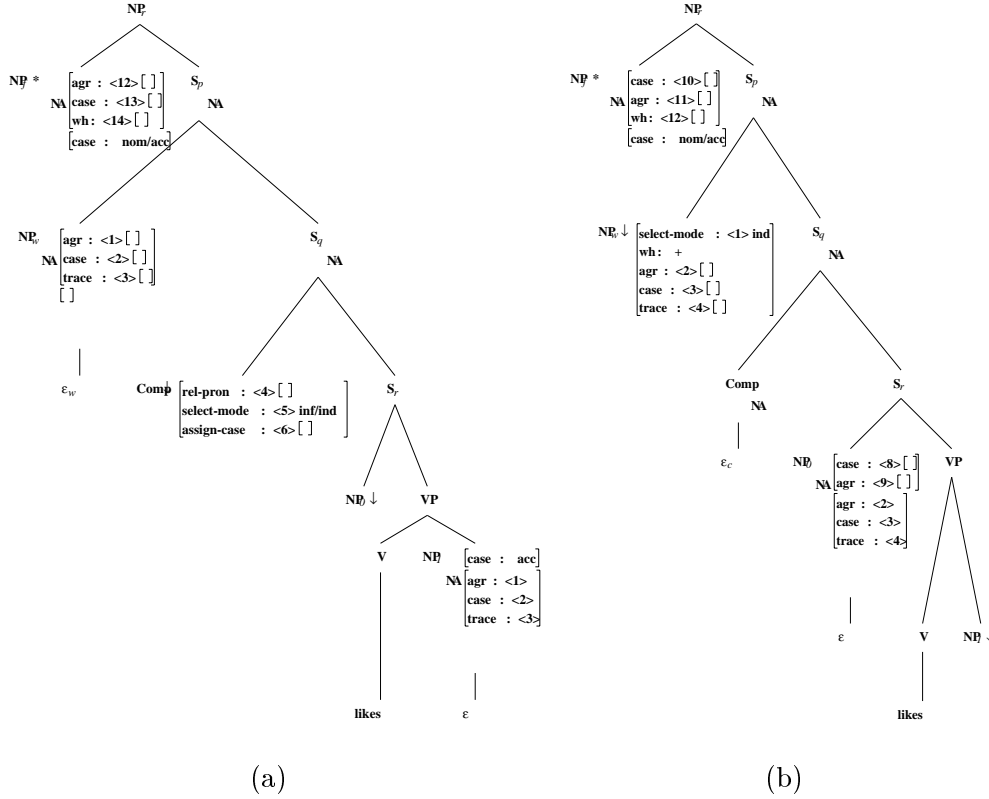


Figure 14.1: Relative clause trees in the transitive tree family: $\beta Nc1nx0Vnx1$ (a) and $\beta N0nx0Vnx1$ (b)

The above analysis is essentially identical to the GB analysis of relative clauses. One aspect of its implementation is that an covert $+ <wh>$ NP and a covert Comp have to be introduced. See (159) and (160) for example.

(159) export exhibitions $[[NP_w \epsilon]_i [\text{that} [\epsilon_i \text{ included high-tech items}]]]$

(160) the export exhibition $[[NP_w \epsilon]_i [\epsilon_C [\text{Muriel planned } \epsilon_i]]]$

The lexicalized nature of XTAG makes it problematic to have trees headed by null strings. Of the two null trees, NP_w and Comp, that we could postulate, the former is definitely more undesirable because it would lead to massive overgeneration, as can be seen in (161) and (162).

(161) $* [NP_w \epsilon]$ did John eat the apple? (as a *wh*-question)

(162) $* \text{I wonder } [[NP_w \epsilon] \text{ Mary likes John}]$ (as an indirect question)

The presence of an initial headed by a null Comp does not lead to problems of overgeneration because relative clauses are the only environment with a Comp substitution node.³

³Complementizers in clausal complementation are introduced by adjunction. See section 8.4.

Consequently, our treatment of relative clauses has different trees to handle relative clauses with an overt extracted *wh*-NP and relative clauses with a covert extracted *wh*-NP. Relative clauses with an overt extracted *wh*-NP involve substitution of a +<**wh**> NP into the NP_w node⁴ and have a Comp node headed by ϵ_C built in. Relative clauses with a covert extracted *wh*-NP have a NP_w node headed by ϵ_w built in and involve substitution into the Comp node. The Comp node that is introduced by substitution can be the ϵ_C (null complementizer), *that*, and *for*.

For example, the tree shown in Figure 14.1(b) is used for the relative clauses shown in sentences (163)-(164), while the tree shown in Figure 14.1(a) is used for the relative clauses in sentences (165)-(168).

(163) the man who Muriel likes

(164) the man whose mother Muriel likes

(165) the man Muriel likes

(166) the book for Muriel to read

(167) the man that Muriel likes

(168) the book Muriel is reading

Cases of PP pied-piping (cf. 169) are handled in a similar fashion by building in a PP_w node.

(169) the demon by whom Muriel was chased

See the tree in Figure 14.2.

14.1 Complementizers and clauses

The co-occurrence constraints that exist between various Comps and the clause type of the clause they occur with are implemented through combinations of different clause types using the <**mode**> feature, the <**select-mode**> feature, and the <**rel-pron**> feature.

Clauses are specified for the <**mode**> feature which indicates the clause type of that clause. Possible values for the <**mode**> feature are **ind**, **inf**, **ppart**, **ger** etc.

Comps are lexically specified for a feature named <**select-mode**>. In addition, the <**select-mode**> feature of the Comp is equated with the <**mode**> feature of its complement S by the following equation:

$$\mathbf{S}_r.\mathbf{t}:\langle\mathbf{mode}\rangle = \mathbf{Comp.t}:\langle\mathbf{select-mode}\rangle$$

The lexical specifications of the Comps are shown below:

- ϵ_C , **Comp.t:**<**select-mode**> = ind/inf/ger/ppart

⁴The feature equation used is $\mathbf{NP}_w.\mathbf{t}:\langle\mathbf{wh}\rangle = +$. Examples of NPs that could substitute under NP_w are *whose mother*, *who*, *whom*, and also *which* but not *when* and *where* which are treated as exhaustive +*wh* PPs.

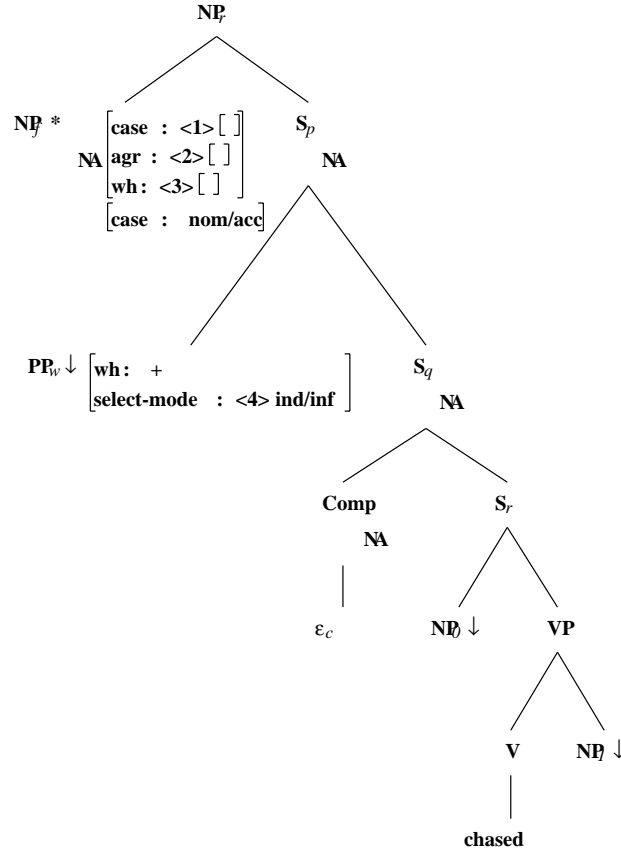


Figure 14.2: Adjunct relative clause tree with PP-pied-piping in the transitive tree family:
 $\beta N_{pxnx0V_{nx1}}$

- *that*, **Comp.t:** $\langle \text{select-mode} \rangle = \text{ind}$
- *for*, **Comp.t:** $\langle \text{select-mode} \rangle = \text{inf}$

The following examples display the co-occurrence constraints which the $\langle \text{select-mode} \rangle$ specifications assigned above implement.

For ϵ_c :

- (170) the book Muriel likes (**S.t:** $\langle \text{mode} \rangle = \text{ind}$)
- (171) a book to like (**S.t:** $\langle \text{mode} \rangle = \text{inf}$)
- (172) the girl reading the book (**S.t:** $\langle \text{mode} \rangle = \text{ger}$)
- (173) the book read by Muriel (**S.t:** $\langle \text{mode} \rangle = \text{ppart}$)

For *for*:

- (174) *the book for Muriel likes (**S.t:**<**mode**>= **ind**)
 (175) a book for Mary to like (**S.t:**<**mode**>= **inf**)
 (176) *the girl for reading the book (**S.t:**<**mode**>= **ger**)
 (177) *the book for read by Muriel (**S.t:**<**mode**>= **ppart**)

For *that*:

- (178) the book that Muriel likes (**S.t:**<**mode**>= **ind**)
 (179) *a book that (Muriel) to like (**S.t:**<**mode**>= **inf**)
 (180) *the girl that reading the book (**S.t:**<**mode**>= **ger**)
 (181) *the book that read by Muriel (**S.t:**<**mode**>= **ppart**)

Relative clause trees that have substitution of \mathbf{NP}_w have the following feature equations:
 $\mathbf{S}_r.\mathbf{t}:\langle\mathbf{mode}\rangle = \mathbf{NP}_w.\mathbf{t}:\langle\mathbf{select-mode}\rangle$
 $\mathbf{NP}_w.\mathbf{t}:\langle\mathbf{select-mode}\rangle = \mathbf{ind}$

The examples that follow are intended to provide the rationale for the above setting of features.

- (182) the boy whose mother chased the cat (**S_r.t:**<**mode**>=**ind**)
 (183) *the boy whose mother to chase the cat (**S_r.t:**<**mode**>=**inf**)
 (184) *the boy whose mother eaten the cake (**S_r.t:**<**mode**>=**ppart**)
 (185) *the boy whose mother chasing the cat (**S_r.t:**<**mode**>= **ger**)
 (186) the boy [whose mother]_i Bill believes ϵ_i to chase the cat
 (**S_r.t:** <**mode**>=**ind**)

The feature equations that appear in trees which have substitution of \mathbf{PP}_w are:
 $\mathbf{S}_r.\mathbf{t}:\langle\mathbf{mode}\rangle = \mathbf{PP}_w.\mathbf{t}:\langle\mathbf{select-mode}\rangle$
 $\mathbf{PP}_w.\mathbf{t}:\langle\mathbf{mode}\rangle = \mathbf{ind}/\mathbf{inf}$ ⁵

Examples that justify the above feature setting follow.

- (187) the person [by whom] this machine was invented (**S_r.t:**<**mode**>=**ind**)
 (188) a baker [in whom]_i PRO to trust ϵ_i (**S_r.t:**<**mode**>= **inf**)
 (189) *the fork [with which] (Geoffrey) eaten the pudding (**S_r.t:**< **mode**>=**ppart**)
 (190) *the person [by whom] (this machine) inventing (**S_r.t:**<**mode**>=**ger**)

⁵As is the case for \mathbf{NP}_w substitution, any +**wh**-PP can substitute under \mathbf{PP}_w . This is implemented by the following equation:

$$\mathbf{PP}_w.\mathbf{t}:\langle\mathbf{wh}\rangle = +$$

Not all cases of pied-piping involve substitution of \mathbf{PP}_w . In some cases, the P may be built in. In cases where part of the pied-piped PP is part of the anchor, it continues to function as an anchor even after pied-piping i.e. the P node and the \mathbf{NP}_w nodes are represented separately.

14.1.1 Further constraints on the null Comp ϵ_C

There are additional constraints on where the null Comp ϵ_C can occur. The null Comp is not permitted in cases of subject extraction unless there is an intervening clause or the relative clause is a reduced relative (**mode** = **ppart/ger**). This can be seen in (191-194).

(191) *the toy [ϵ_i [ϵ_C [ϵ_i likes Dafna]]]

(192) the toy [ϵ_i [ϵ_C Fred thinks [ϵ_i likes Dafna]]]

(193) the boy [ϵ_i [ϵ_C [ϵ_i eating the guava]]]

(194) the guava [ϵ_i [ϵ_C [ϵ_i eaten by the boy]]]

To model this paradigm, the feature **<rel-pron>** is used in conjunction with the following equations:

- **S_r.t:<rel-pron>** = **Comp.t:<rel-pron>**
- **S_r.b:<rel-pron>** = **S_r.b:<mode>**
- **Comp.b:<rel-pron>** = **ppart/ger/adj-clause** (for ϵ_C)

The full set of the equations shown above is only present in Comp substitution trees involving subject extraction. So (195) will not be ruled out.

(195) the toy [ϵ_i [ϵ_C [ϵ_i Dafna likes ϵ_i]]]

The feature mismatch induced by the above equations is not remedied by adjunction of just any S-adjunct because all other S-adjuncts are transparent to the **<rel-pron>** feature because of the following equation:

S_m.b:<rel-pron> = **S_f.t:<rel-pron>**

14.2 Reduced Relatives

Reduced relatives are permitted only in cases of subject-extraction. Past participial reduced relatives are only permitted on passive clauses. See (196-203).

(196) the toy [ϵ_i [ϵ_C [ϵ_i playing the banjo]]]

(197) *the instrument [ϵ_i [ϵ_C [ϵ_i Amis playing ϵ_i]]]

(198) *the day [ϵ_w [ϵ_C [ϵ_i Amis playing the banjo]]]

(199) the apple [ϵ_i [ϵ_C [ϵ_i eaten by Dafna]]]

(200) *the child [ϵ_i [ϵ_C [ϵ_i the apple eaten by ϵ_i]]]

(201) *the day [ϵ_w [ϵ_C [ϵ_i Amis eaten the apple]]]

(202) *the apple [ϵ_i [ϵ_C [Dafna eaten ϵ_i]]]

(203) *the child [ϵ_i [ϵ_C [ϵ_i eaten the apple]]]

These restrictions are built into the $\langle \mathbf{mode} \rangle$ specifications of **S.t**. So non-passive cases of subject extraction have the following feature equation:

S_r.t: $\langle \mathbf{mode} \rangle = \mathbf{ind/ger/inf}$

Passive cases of subject extraction have the following feature equation:

S_r.t: $\langle \mathbf{mode} \rangle = \mathbf{ind/ger/ppart/inf}$

Finally, all cases of non-subject extraction have the following feature equation:

S_r.t: $\langle \mathbf{mode} \rangle = \mathbf{ind/inf}$

14.2.1 Restrictive vs. Non-restrictive relatives

The English XTAG grammar does not contain any syntactic distinction between restrictive and non-restrictive relatives because we believe this to be a semantic and/or pragmatic difference.

14.3 External syntax

A relative clause can combine with the NP it modifies in at least the following two ways:

(204) [the [toy [ϵ_i [ϵ_C [Dafna likes ϵ_i]]]]]

(205) [[the toy] [ϵ_i [ϵ_C [Dafna likes ϵ_i]]]]

Based on cases like (206) and (207), which are problematic for the structure in (204), the structure in (205) is adopted.

(206) [[the man and the woman] [who met on the bus]]

(207) [[the man and the woman] [who like each other]]

As it stands, the RC analysis sketched so far will combine in two ways with the Determiner tree shown in Figure (207),⁶ giving us both the possibilities shown in (204) and (205). In order to block the structure exemplified in (204), the feature $\langle \mathbf{rel-clause} \rangle$ is used in combination with the following equations.

On the RC:

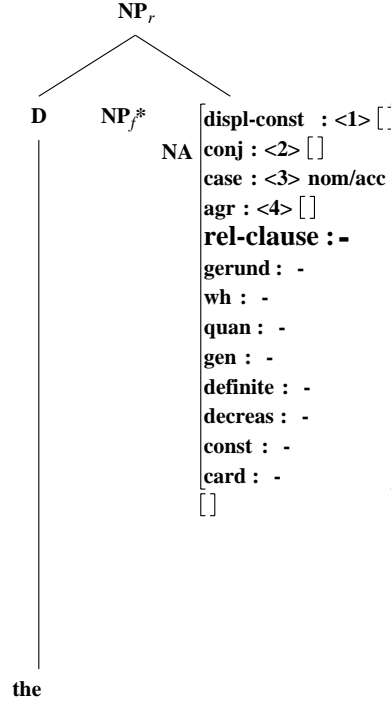
NP_r.b: $\langle \mathbf{rel-clause} \rangle = +$

On the Determiner tree:

NP_f.t: $\langle \mathbf{rel-clause} \rangle = -$

Together, these equations block introduction of the determiner above the relative clause.

⁶The determiner tree shown has the $\langle \mathbf{rel-clause} \rangle$ feature built in. The RC analysis would give two parses in the absence of this feature.

Figure 14.3: Determiner tree with $\langle \text{rel-clause} \rangle$ feature: βDnx

14.4 Other Issues

14.4.1 Interaction with adjoined Comps

The XTAG analysis now has two different ways of introducing a complementizer like *that* or *for*, depending upon whether it occurs in a relative clause or in sentential complementation. Relative clause complementizers substitute in (using the tree αComp), while sentential complementizers adjoin in (using the tree βCOMPs). Cases like (208) where both kinds of complementizers illicitly occur together are blocked.

(208) *the book [ϵ_{w_i} [that [that [Muriel wrote ϵ_i]]]]

This is accomplished by setting the $\mathbf{S_r.t:comp}$ feature in the relative clause tree to **nil**. The $\mathbf{S_r.t:comp}$ feature of the auxiliary tree that introduces (the sentential complementation) *that* is set to **that**. This leads to a feature clash ruling out (208). On the other hand, if a sentential complement taking verb is adjoined in at S_r , this feature clash goes away (cf. 209).

(209) the book [ϵ_{w_i} [that Beth thinks [that [Muriel wrote ϵ_i]]]]

14.4.2 Adjunction on PRO

Adjunction on PRO, which would yield the ungrammatical (210) is blocked.

(210) *I want [[PRO [who Muriel likes] to read a book]].

This is done by specifying the $\langle \text{case} \rangle$ feature of \mathbf{NP}_f to be **nom/acc**. The $\langle \text{case} \rangle$ feature of PRO is **null**. This leads to a feature clash and blocks adjunction of relative clauses on to PRO.

14.4.3 Adjunct relative clauses

Two types of trees to handle adjunct relative clauses exist in the XTAG grammar: one in which there is \mathbf{PP}_w substitution with a null **Comp** built in and one in which there is a null \mathbf{NP}_w built in and a **Comp** substitutes in. There is no \mathbf{NP}_w substitution tree with a null **Comp** built in. This is because of the contrast between (211) and (212).

(211) the day [[on whose predecessor] [ϵ_C [Muriel left]]]

(212) *the day [[whose predecessor] [ϵ_C [Muriel left]]]

In general, adjunct relatives are not possible with an overt \mathbf{NP}_w . We do not consider (213) and (214) to be counterexamples to the above statements because we consider *where* and *when* to be exhaustive **PP**s that head a **PP** initial tree.

(213) the place [where [ϵ_C [Muriel wrote her first book]]]

(214) the time [when [ϵ_C [Muriel lived in Bryn Mawr]]]

14.4.4 ECM

Cases where *for* assigns exceptional case (cf. 215, 216) are handled.

(215) a book [ϵ_{w_i} [for [Muriel to read ϵ_i]]]

(216) the time [ϵ_{w_i} [for [Muriel to leave Haverford]]]

The assignment of case by *for* is implemented by a combination of the following equations.

$\mathbf{Comp.t}:\langle \text{assign-case} \rangle = \text{acc}$

$\mathbf{S_r.t}:\langle \text{assign-case} \rangle = \mathbf{Comp.t}:\langle \text{assign-case} \rangle$

$\mathbf{S_r.b}:\langle \text{assign-case} \rangle = \mathbf{NP_0.t}:\langle \text{case} \rangle$

14.5 Cases not handled

14.5.1 Partial treatment of free-relatives

Free relatives are only partially handled. All free relatives on non-subject positions and some free relatives on subject positions are handled. The structure assigned to free relatives treats the extracted *wh*-NP as the head NP of the relative clause. The remaining relative clause modifies this extracted *wh*-NP (cf. 217-219).

(217) what(ever) [ϵ_{w_i} [ϵ_C [Mary likes ϵ_i]]]

(218) where(ever) [ϵ_w [ϵ_C [Mary lives]]]

(219) who(ever) [ϵ_{w_i} [ϵ_C [Muriel thinks [ϵ_i likes Mary]]]]

However, simple subject extractions without further embedding are not handled (cf. 220).

(220) who(ever) [ϵ_{w_i} [ϵ_C [ϵ_i likes Bill]]]

This is because (219) is treated exactly like the ungrammatical (221).

(221) *the person [ϵ_{w_i} [ϵ_C [ϵ_i likes Bill]]]

14.5.2 Adjunct P-stranding

The following cases of adjunct preposition stranding are not handled (cf. 222, 223).

(222) the pen Muriel wrote this letter with

(223) the street Muriel lives on

Adjuncts are not built into elementary trees in XTAG. So there is no clean way to represent adjunct preposition stranding. A better solution is, probably, available if we make use of multi-component adjunction.

14.5.3 Overgeneration

The following ungrammatical sentences are currently being accepted by the XTAG grammar. This is because no clean and conceptually attractive way of ruling them out is obvious to us.

14.5.3.1 *how* as *wh*-NP

In standard American English, *how* is not acceptable as a relative pronoun (cf. 224).

(224) *the way [how [ϵ_C [PRO to solve this problem]]]

However, (224) is accepted by the current grammar. The only way to rule (224) out would be to introduce a special feature devoted to this purpose. This is unappealing. Further, there exist speech registers/dialects of English, where (224) is acceptable.

14.5.3.2 *for*-trace effects

(225) is ungrammatical, being an instance of a violation of the *for*-trace filter of early transformational grammar.

(225) the person [ϵ_{w_i} [for [ϵ_i to read the book]]]

The XTAG grammar currently accepts (225).⁷

⁷It may be of some interest that (225) is acceptable in certain dialects of Belfast English.

14.5.3.3 Internal head constraint

Relative clauses in English (and in an overwhelming number of languages) obey a ‘no internal head’ constraint. This constraint is exemplified in the contrast between (226) and (227).

(226) the person [who_i [ϵ_C Muriel likes ϵ_i]]

(227) *the person [[which person] $_i$ [ϵ_C Muriel likes ϵ_i]]

We know of no good way to rule (227) out, while still ruling (228) in.

(228) the person [[whose mother] $_i$ [ϵ_C Muriel likes ϵ_i]]

Dayal (1996) suggests that ‘full’ NPs such as *which person* and *whose mother* are R-expressions while *who* and *whose* are pronouns. R-expressions, unlike pronouns, are subject to Condition C. (226) is, then, ruled out as a violation of Condition C since *the person* and *which person* are co-indexed and *the person* c-commands *which person*. If we accept Dayal’s argument, we have a principled reason for allowing overgeneration of relative clauses that violate the internal head constraint, the reason being that the XTAG grammar does generate binding theory violations.

14.5.3.4 Overt Comp constraint on stacked relatives

Stacked relatives of the kind in (229) are handled.

(229) [[the book [that Bill likes]] [which Mary wrote]]

There is a constraint on stacked relatives: all but the relative clause closest to the head-NP must have either an overt **Comp** or an overt **NP_w**. Thus (230) is ungrammatical.

(230) *[[the book [that Bill likes]] [Mary wrote]]

Again, no good way of handling this constraint is known to us currently.

Chapter 15

Adjunct Clauses

Adjunct clauses include subordinate clauses (i.e. those with overt subordinating conjunctions), purpose clauses and participial adjuncts.

Subordinating conjunctions each select four trees, allowing them to appear in four different positions relative to the matrix clause. The positions are (1) before the matrix clause, (2) after the matrix clause, (3) before the VP, surrounded by two punctuation marks, and (4) after the matrix clause, separated by a punctuation mark. Each of these trees is shown in Figure 15.1.

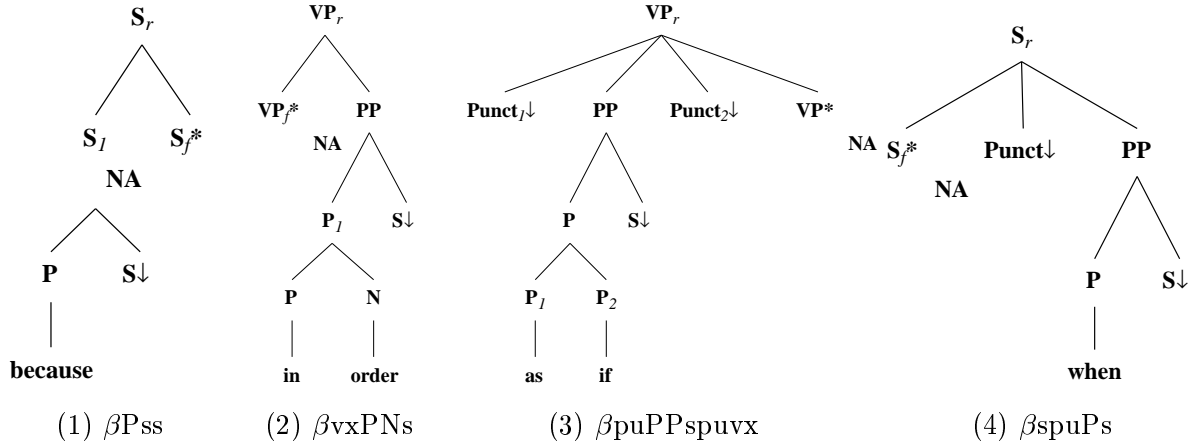


Figure 15.1: Auxiliary Trees for Subordinating Conjunctions

Sentence-initial adjuncts adjoin at the root S of the matrix clause, while sentence-final adjuncts adjoin at a VP node. In this, the XTAG analysis follows the findings on the attachment sites of adjunct clauses for conditional clauses ([Iatridou, 1991]) and for infinitival clauses ([Browning, 1987]). One compelling argument is based on Binding Condition C effects. As can be seen from examples (231)-(233) below, no Binding Condition violation occurs when the adjunct is sentence initial, but the subject of the matrix clause clearly governs the adjunct clause when it is in sentence final position and co-indexation of the pronoun with the subject of the adjunct clause is impossible.

(231) Unless she_i hurries, Mary_i will be late for the meeting.

(232) *She_i will be late for the meeting unless Mary_i hurries.

(233) Mary_i will be late for the meeting unless she_i hurries.

We had previously treated subordinating conjunctions as a subclass of *conjunction*, but are now assigning them the POS *preposition*, as there is such clear overlap between words that function as prepositions (taking NP complements) and subordinating conjunctions (taking clausal complements). While there are some prepositions which only take NP complements and some which only take clausal complements, many take both as shown in examples (234)-(237), and it seems to be artificial to assign them two different parts-of-speech.

(234) Helen left before the party.

(235) Helen left before the party began.

(236) Since the election, Bill has been elated.

(237) Since winning the election, Bill has been elated.

Each subordinating conjunction selects the values of the **<mode>** and **<comp>** features of the subordinated S. The **<mode>** value constrains the types of clauses the subordinating conjunction may appear with and the **<comp>** value constrains the complementizers which may adjoin to that clause. For instance, indicative subordinate clauses may appear with the complementizer *that* as in (238), while participial clauses may not have any complementizers (239).

(238) Midge left that car so that Sam could drive to work.

(239) *Since that seeing the new VW, Midge could think of nothing else.

15.0.4 Multi-word Subordinating Conjunctions

We extracted a list of multi-word conjunctions, such as *as if*, *in order*, and *for all (that)*, from [Quirk *et al.*, 1985]. For the most part, the components of the complex are all anchors, as shown in Figures 15.2(a). In one case, *as ADV as*, there is a great deal of latitude in the choice of adverb, so this is a substitution site (Figures 15.2(b)). This multi-anchor treatment is very similar to that proposed for idioms in [Abeillé and Schabes, 1989], and the analysis of light verbs in the XTAG grammar (see section 6.15).

15.1 “Bare” Adjunct Clauses

“Bare” adjunct clauses do not have an overt subordinating conjunction, but are typically parallel in meaning to clauses with subordinating conjunctions. For this reason, we have elected to handle them using the same trees shown above, but with null anchors. They are selected at the same time and in the same way the *PRO* tree is, as they all have *PRO* subjects. Three values of **<mode>** are licensed: **inf** (infinitive), **ger** (gerundive) and **ppart** (past participial).¹ They interact with complementizers as follows:

¹We considered allowing bare indicative clauses, such as *He died that others may live*, but these were considered too archaic to be worth the additional ambiguity they would add to the grammar.

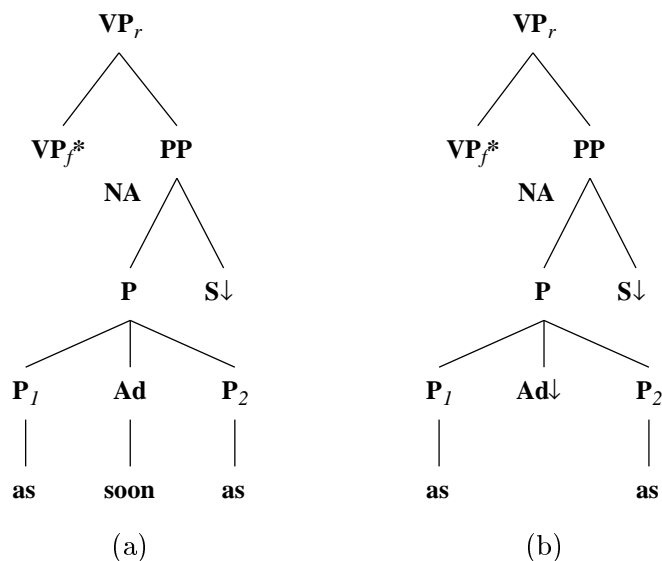


Figure 15.2: Trees Anchored by Subordinating Conjunctions: $\beta vxPARBPs$ and $\beta vxParbPs$

- Participial complements do not license any complementizers:²
 - (240) [Destroyed by the fire], the building still stood.
 - (241) The fire raged for days [destroying the building].
 - (242) *[That destroyed by the fire], the building still stood.
- Infinitival adjuncts, including purpose clauses, are licensed both with and without the complementizer *for*.
 - (243) Harriet bought a Mustang [to impress Eugene].
 - (244) [To impress Harriet], Eugene dyed his hair.
 - (245) Traffic stopped [for Harriet to cross the street].

15.2 Discourse Conjunction

The CONJs auxiliary tree is used to handle ‘discourse’ conjunction, as in sentence (246). Only the coordinating conjunctions (*and*, *or* and *but*) are allowed to adjoin to the roots of matrix sentences. Discourse conjunction with *and* is shown in the derived tree in Figure 15.4.

(246) And Truffula trees are what everyone needs! [Seuss, 1971]

²While these sound a bit like extraposed relative clauses (see [Kroch and Joshi, 1987]), those move only to the right and adjoin to S; as these clauses are equally grammatical both sentence-initially and sentence-finally, we are analyzing them as adjunct clauses.

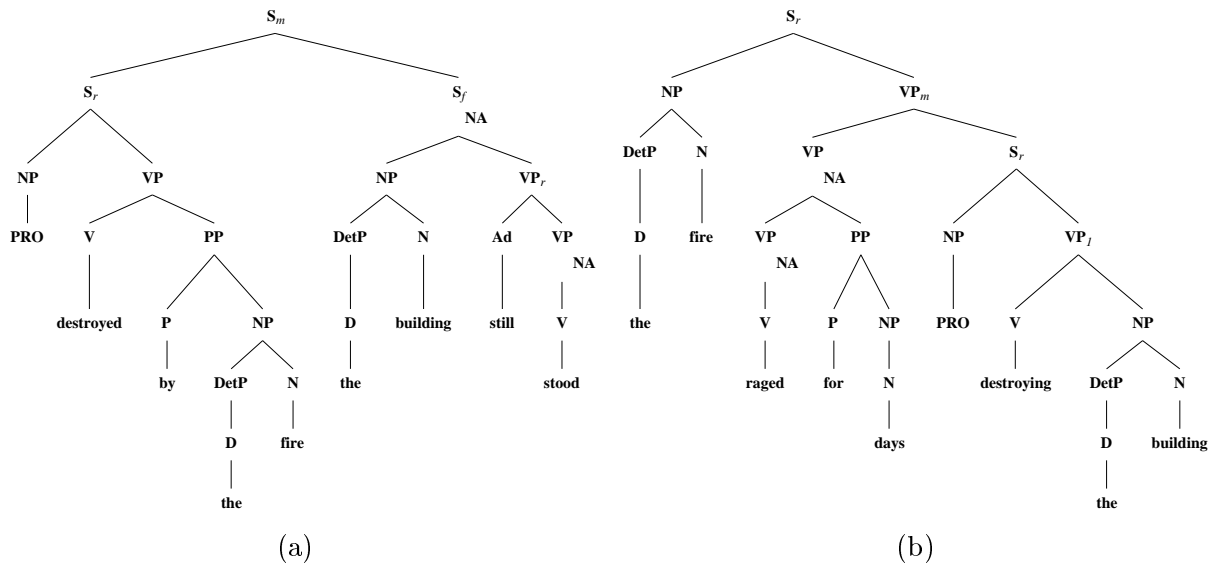


Figure 15.3: Sample Participial Adjuncts

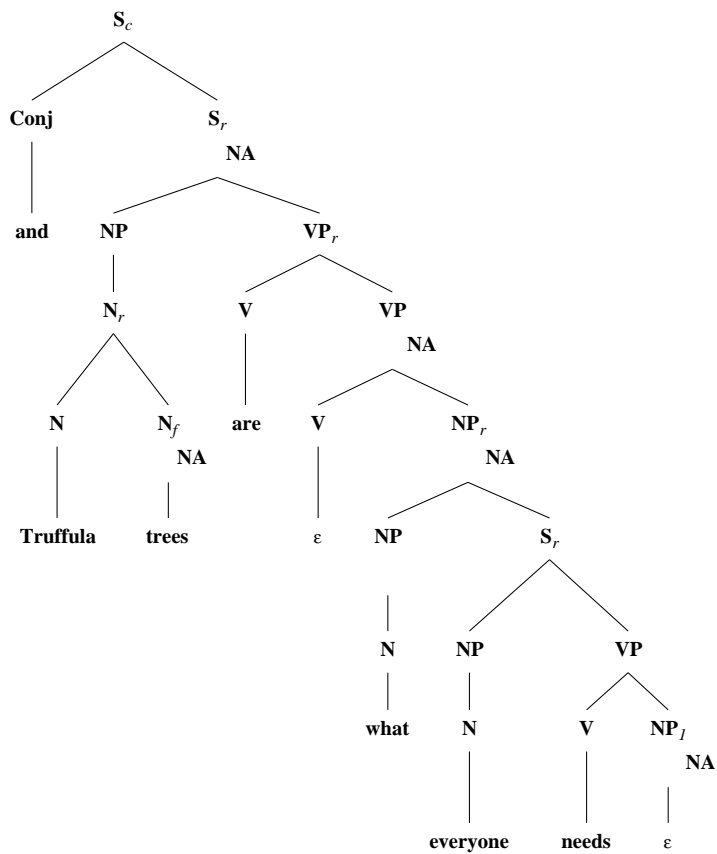


Figure 15.4: Example of discourse conjunction, from Seuss' *The Lorax*

Chapter 16

Imperatives

Imperatives in English do not require overt subjects. The subject in imperatives is second person, i.e. *you*, whether it is overt or not, as is clear from the verbal agreement and the interpretation. Imperatives with overt subjects can be parsed using the trees already needed for declaratives. The imperative cases in which the subject is not overt are handled by the imperative trees discussed in this section.

The imperative trees in English XTAG grammar are identical to the declarative tree except that the NP₀ subject position is filled by an ϵ , the NP₀ <**agr pers**> feature is set in the tree to the value **2nd** and the <**mode**> feature on the root node has the value **imp**. The value for <**agr pers**> is hardwired into the epsilon node and insures the proper verbal agreement for an imperative. The <**mode**> value of **imp** on the root node is recognized as a valid mode for a matrix clause. The **imp** value for <**mode**> also allows imperatives to be blocked from appearing as embedded clauses. Figure 16.1 is the imperative tree for the transitive tree family.

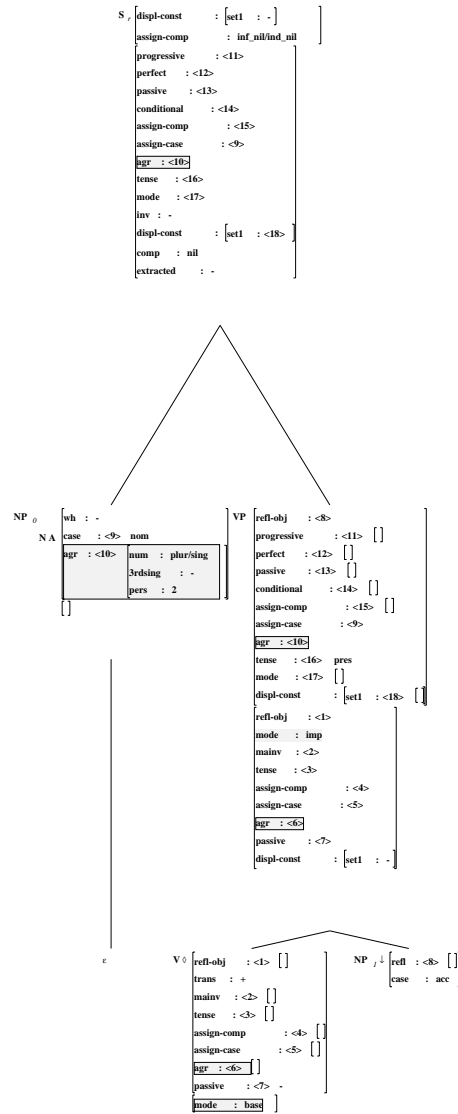


Figure 16.1: Transitive imperative tree: $\alpha\text{Inx0Vnx1}$

Chapter 17

Gerund NP's

There are two types of gerunds identified in the linguistics literature. One is the class of *derived nominalizations* (also called *nominal gerundives* or *action nominalizations*) exemplified in (247), which instantiates the direct object within an *of* PP. The other is the class of so-called *sentential* or *VP gerundives* exemplified in (248). In the English XTAG grammar, the derived nominalizations are termed **determiner gerunds**, and the sentential or VP gerunds are termed **NP gerunds**.

(247) Some think that **the selling of bonds** is beneficial.

(248) Are private markets approving of **Washington bashing Wall Street**?

Both types of gerunds exhibit a similar distribution, appearing in most places where NP's are allowed.¹ The bold face portions of sentences (249)–(251) show examples of gerunds as a subject and as the object of a preposition.

(249) **Avoiding such losses** will take a monumental effort.

(250) **Mr. Nolen's wandering** doesn't make him a weirdo.

(251) Are private markets approving of **Washington bashing Wall Street**?

The motivation for splitting the gerunds into two classes is semantic as well as structural in nature. Semantically, the two gerunds are in sharp contrast with each other. NP gerunds refer to an action, i.e., a way of doing something, whereas determiner gerunds refer to a fact. Structurally, there are a number of properties (extensively discussed in [Lees, 1960]) that show that NP gerunds have the syntax of verbs, whereas determiner gerunds have the syntax of basic nouns. Firstly, the fact that the direct object of the determiner gerund can only appear within an *of* PP suggests that the determiner gerund, like nouns, is not a case assigner and needs insertion of the preposition *of* for assignment of case to the direct object. NP gerunds, like verbs, need no such insertion and can assign case to their direct object. Secondly, like nouns, only determiner gerunds can appear with articles (cf. example (252) and (253)). Thirdly, determiner gerunds, like nouns, can be modified by adjectives (cf. example (254)), whereas

¹an exception being the NP positions in “equative BE” sentences, such as, *John is my father*.

NP gerunds, like verbs, resist such modification (cf. example (255)). Fourthly, nouns, unlike verbs, cannot co-occur with aspect (cf. example (256) and (257)). Finally, only NP gerunds, like verbs, can take adverbial modification (cf. example (258) and (259)).

- | | |
|-----------------------------------------------|------------------------|
| (252) ... the proving of the theorem.... | (det ger with article) |
| (253) * ... the proving the theorem.... | (NP ger with article) |
| (254) John's rapid writing of the book.... | (det ger with Adj) |
| (255) * John's rapid writing the book.... | (NP ger with Adj) |
| (256) * John's having written of the book.... | (det ger with aspect) |
| (257) John having written the book.... | (NP ger with aspect) |
| (258) * His writing of the book rapidly.... | (det ger with Adverb) |
| (259) His writing the book rapidly.... | (NP ger with Adverb) |

In English XTAG, the two types of gerunds are assigned separate trees within each tree family, but in order to capture their similar distributional behavior, both are assigned NP as the category label of their top node. The feature **gerund** = +/- distinguishes gerund NP's from regular NP's where needed.² The determiner gerund and the NP gerund trees are discussed in section (17.1) and (17.2) respectively.

17.1 Determiner Gerunds

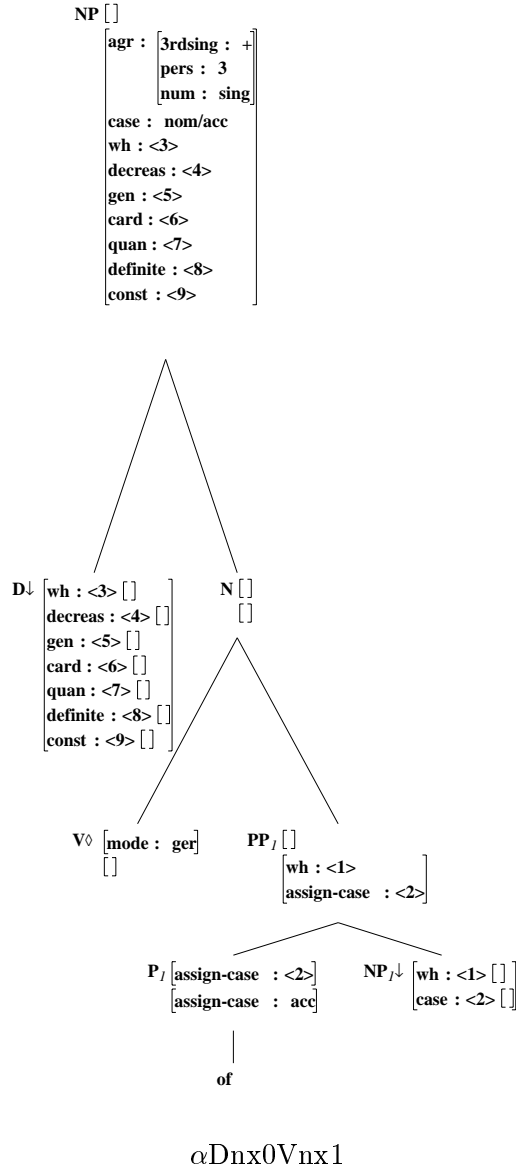
The determiner gerund tree in Figure 17.1 is anchored by a V, capturing the fact that the gerund is derived from a verb. The verb projects an N and instantiates the direct object as an *of* PP. The nominal category projected by the verb can now display all the syntactic properties of basic nouns, as discussed above. For example, it can be straightforwardly modified by adjectives; it cannot co-occur with aspect; and it can appear with articles. The only difference of the determiner gerund nominal with the basic nominals is that the former cannot occur without the determiner, whereas the latter can. The determiner gerund tree therefore has an initial D modifying the N.³ It is used for gerunds such as the ones in bold face in sentences (260), (261) and (262).

The D node can take a simple determiner (cf. example (260)), a genitive pronoun (cf. example (261)), or a genitive NP (cf. example (262)).⁴

²This feature is also needed to restrict the selection of gerunds in NP positions. For example, the subject and object NP's in the "equative BE" tree (Tnx0BEnx1) cannot be filled by gerunds, and are therefore assigned the feature **gerund** = -, which prevents gerunds (which have the feature **gerund** = +) from substituting into these NP positions.

³Note that the determiner can adjoin to the gerund only from *within* the gerund tree. Adjunction of determiners to the gerund root node is prevented by constraining determiners to only select NP's with the feature **gerund** = -. This rules out sentences like *Private markets approved of (*the) [the selling of bonds]*.

⁴The trees for genitive pronouns and genitive NP's have the root node labelled as D because they seem to function as determiners and as such, sequence with the rest of the determiners. See Chapter 18 for discussion on determiner trees.

Figure 17.1: Determiner Gerund tree from the transitive tree family: $\alpha Dnx0Vnx1$

- (260) Some think that **the selling of bonds** is beneficial.
- (261) **His painting of Mona Lisa** is highly acclaimed.
- (262) Are private markets approving of **Washington's bashing of Wall Street**?

17.2 NP Gerunds

NP gerunds show a number of structural peculiarities, the crucial one being that they have the internal properties of sentences. In the English XTAG grammar, we adopt a position similar

to that of [Rosenbaum, 1967] and [Emonds, 1970] – that these gerunds are NP’s exhaustively dominating a clause. Consequently, the tree assigned to the transitive NP gerund tree (cf. Figure 17.2) looks exactly like the declarative transitive tree for clauses except for the root node label and the feature values. The anchoring verb projects a VP. Auxiliary adjunction is allowed, subject to one constraint – that the highest verb in the verbal sequence be in gerundive form, regardless of whether it is a main or auxiliary verb. This constraint is implemented by requiring the topmost VP node to be **<mode> = ger**. In the absence of any adjunction, the anchoring verb itself is forced to be gerundive. But if the verbal sequence has more than one verb, then the sequence and form of the verbs is limited by the restrictions that each verb in the sequence imposes on the succeeding verb. The nature of these restrictions for sentential clauses, and the manner in which they are implemented in XTAG, are both discussed in Chapter 20. The analysis and implementation discussed there differs from that required for gerunds only in one respect – that the highest verb in the verbal sequence is required to be **<mode> = ger**.

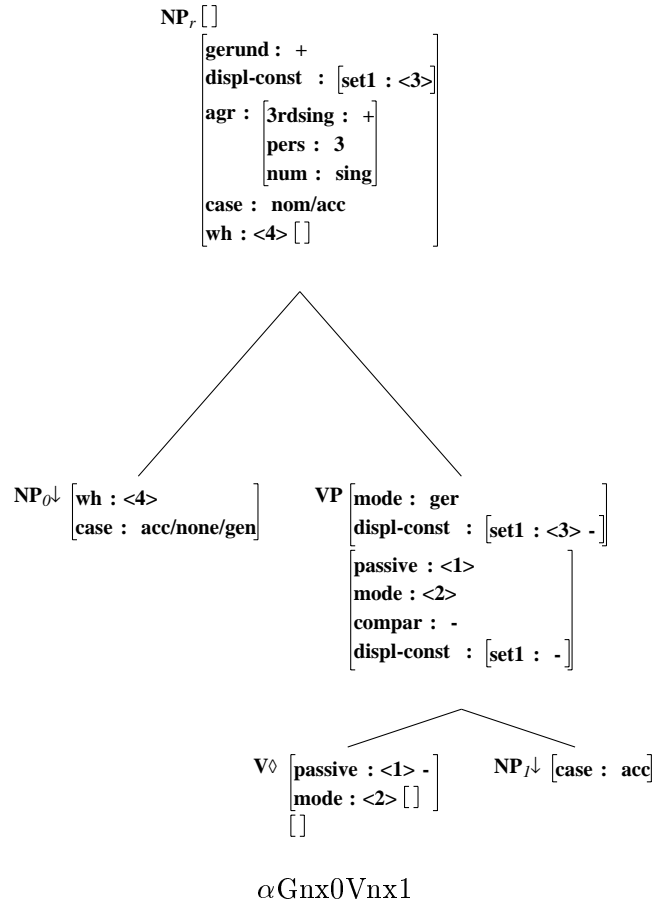


Figure 17.2: NP Gerund tree from the transitive tree family: $\alpha\text{Gnx0Vnx1}$

Additionally, the subject in the NP gerund tree is required to have **<case>=acc/none/gen**, i.e., it can be either a PRO (cf. example 263), a genitive NP (cf. example 264), or an accusative NP (cf. example 265). The whole NP formed by the gerund can occur in either nominative or accusative positions.

- (263) ... John does not like **wearing a hat**.
- (264) Are private markets approving of **Washington's bashing Wall Street**?
- (265) Mother disapproved of **me wearing such casual clothes**.

One question that arises with respect to gerunds is whether there is anything special about their distribution as compared to other types of NP's. In fact, it appears that gerund NP's can occur in any NP position. Some verbs might not seem to be very accepting of gerund NP arguments, as in (266) below, but we believe this to be a semantic incompatibility rather than a syntactic problem since the same structures are possible with other lexical items.

- (266) ? [_{NP}John's tinkering_{NP}] ran.
- (267) [_{NP}John's tinkering_{NP}] worked.

By having the root node of gerund trees be NP, the gerunds have the same distribution as any other NP in the English XTAG grammar without doing anything exceptional. The clause structure is captured by the form of the trees and by inclusion in the tree families.

17.3 Gerund Passives

It was mentioned above that the NP gerunds display certain clausal properties. It is therefore not surprising that they too have their own set of transformationally related structures. For example, NP gerunds allow passivization just like their sentential counterparts (cf. examples (268) and (269)).

- (268) The lawyers objected to **the slanderous book being written by John**.
- (269) Susan could not forget **having been embarrassed by the vicar**.

In the English XTAG grammar, gerund passives are treated in an almost exactly similar fashion to sentential passives, and are assigned separate trees within the appropriate tree families. The passives occur in pairs, one with the *by* phrase, and another without it. There are two feature restrictions imposed on the passive trees: (a) only verbs with **<mode> = ppart** (i.e., verbs with passive morphology) can be the anchors, and (b) the highest verb in the verb sequence is required to be **<mode> = ger**. The two restrictions, together, ensure the selection of only those sequences of auxiliary verb(s) that select **<mode> = ppart** and **<passive> = +** (such as *being* or *having been* but NOT *having*). The passive trees are assumed to be related to only the NP gerund trees (and not the determiner gerund trees), since passive structures involve movement of some object to the subject position (in a movement analysis). Also, like the sentential passives, gerund passives are found in most tree families that have a direct object in the declarative tree. Figure 17.3 shows the gerund passive trees for the transitive tree family.

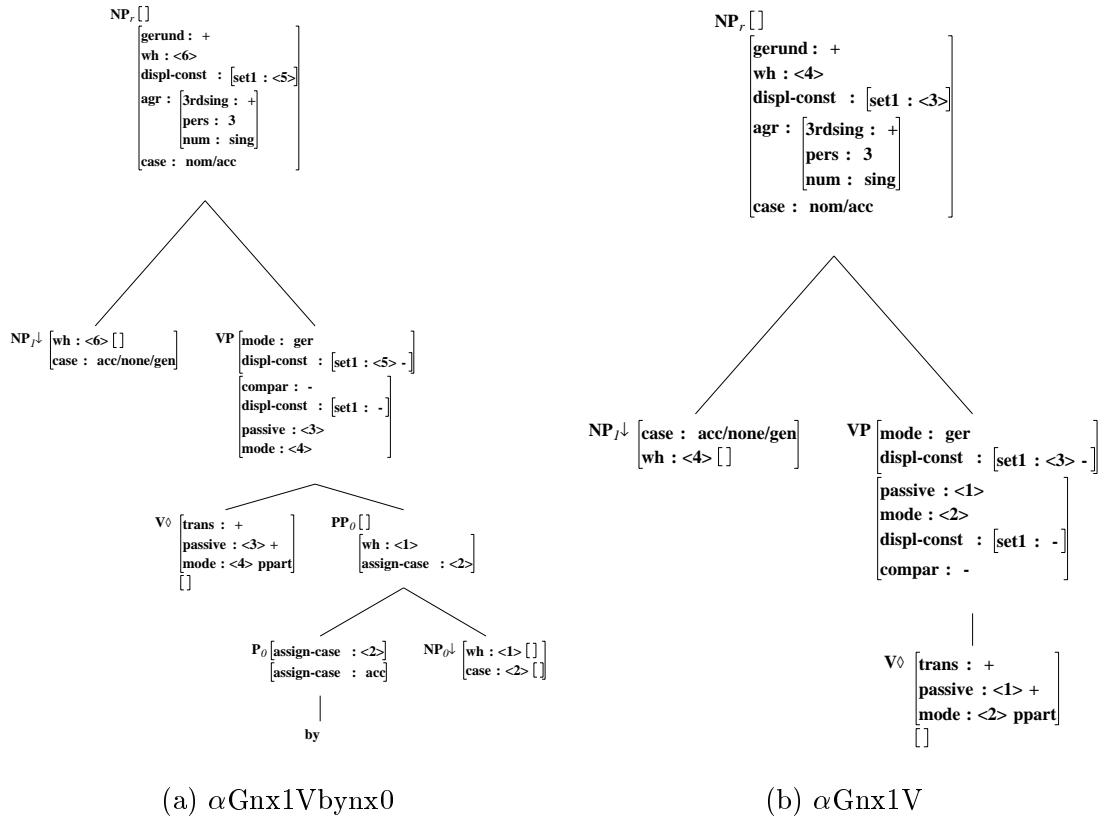


Figure 17.3: Passive Gerund trees from the transitive tree family: $\alpha\text{Gnx1Vbyn0}$ (a) and αGnx1V (b)

Part IV

Other Constructions

Chapter 18

Determiners and Noun Phrases

In our English XTAG grammar,¹ all nouns select the noun phrase (NP) tree structure shown in Figure 18.1. Common nouns do not require determiners in order to form grammatical NPs. Rather than being ungrammatical, singular countable nouns without determiners are restricted in interpretation and can only be interpreted as mass nouns. Allowing all nouns to head determinerless NPs correctly treats the individuation in countable NPs as a property of determiners. Common nouns have negative(“-”) values for determiner features in the lexicon in our analysis and can only acquire a positive(“+”) value for those features if determiners adjoin to them. Other types of NPs such as pronouns and proper nouns have been argued by Abney [Abney, 1987] to either be determiners or to move to the determiner position because they exhibit determiner-like behavior. We can capture this insight in our system by giving pronouns and proper nouns positive values for determiner features. For example pronouns and proper nouns would be marked as definite, a value that NPs containing common nouns can only obtain by having a definite determiner adjoin. In addition to the determiner features, nouns also have values for features such as reflexive (**refl**), case, pronoun (**pron**) and conjunction (**conj**).

A single tree structure is selected by simple determiners, an auxiliary tree which adjoins to NP. An example of this determiner tree anchored by the determiner *these* is shown in Figure 18.2. In addition to the determiner features the tree in Figure 18.2 has noun features such as **case** (see section 4.4.2), the **conj** feature to control conjunction (see Chapter 21), **rel-clause**— (see Chapter 14) and **gerund**— (see Chapter 17) which prevent determiners from adjoining on top of relative clauses and gerund NPs respectively, and the **displ-const** feature which is used to simulate multi-component adjunction.

Complex determiners such as genitives and partitives also anchor tree structures that adjoin to NP. They differ from the simple determiners in their internal complexity. Details of our treatment of these more complex constructions appear in Sections 18.3 and 18.4. Sequences of determiners, as in the NPs *all her dogs* or *those five dogs* are derived by multiple adjunctions of the determiner tree, with each tree anchored by one of the determiners in the sequence. The order in which the determiner trees can adjoin is controlled by features.

This treatment of determiners as adjoining onto NPs is similar to that of [Abeillé, 1990], and allows us to capture one of the insights of the DP hypothesis, namely that determiners select NPs as complements. In Figure 18.2 the determiner and its NP complement appear in

¹A more detailed discussion of this analysis can be found in [Hockey and Mateyak, 1998].

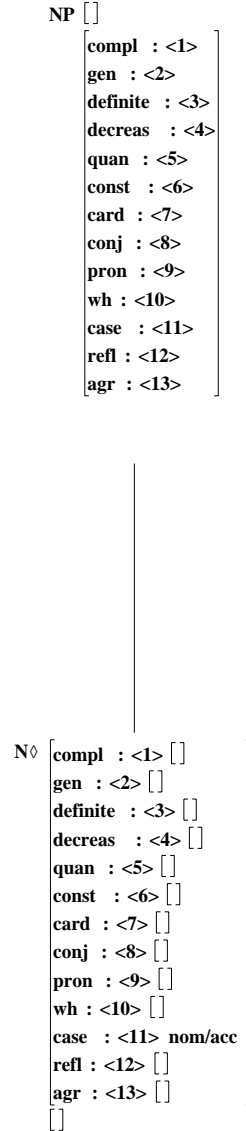


Figure 18.1: NP Tree

the configuration that is typically used in LTAG to represent selectional relationships. That is, the head serves as the anchor of the tree and its complement is a sister node in the same elementary tree.

The XTAG treatment of determiners uses nine features for representing their properties: definiteness (**definite**), quantity (**quan**), cardinality (**card**), genitive (**gen**), decreasing (**decreas**), constancy (**const**), **wh**, agreement (**agr**), and complement (**compl**). Seven of these features were developed by semanticists for their accounts of semantic phenomena ([Keenan and Stavi, 1986], [Barwise and Cooper, 1981], [Partee *et al.*, 1990]), another was developed for

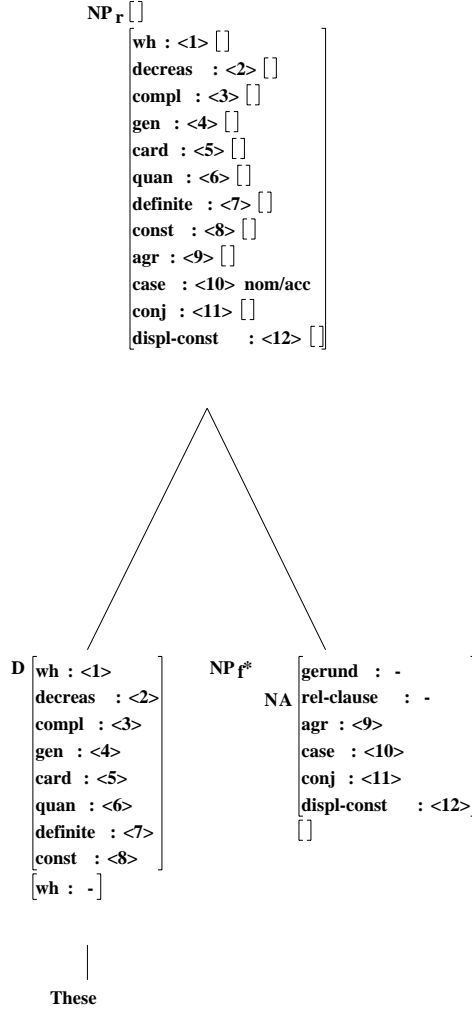


Figure 18.2: Determiner Trees with Features

a semantic account of determiner negation by one of the authors of this determiner analysis ([Mateyak, 1997]), and the last is the familiar agreement feature. When used together these features also account for a substantial portion of the complex patterns of English determiner sequencing. Although we do not claim to have exhaustively covered the sequencing of determiners in English, we do cover a large subset, both in terms of the phenomena handled and in terms of corpus coverage. The XTAG grammar has also been extended to include complex determiner constructions such as genitives and partitives using these determiner features.

Each determiner carries with it a set of values for these features that represents its own properties, and a set of values for the properties of NPs to which can adjoin. The features are crucial to ordering determiners correctly. The semantic definitions underlying the features are given below.

Definiteness: Possible Values [+/-].

A function f is definite iff f is non-trivial and whenever $f(s) \neq \emptyset$ then it is always the intersection of one or more individuals. [Keenan and Stavi, 1986]

Quantity: Possible Values [+/-].

If A and B are sets denoting an NP and associated predicate, respectively; E is a domain in a model M , and F is a bijection from M_1 to M_2 , then we say that a determiner satisfies the constraint of quantity if $\text{Det}_{E_1}AB \leftrightarrow \text{Det}_{E_2}F(A)F(B)$. [Partee *et al.*, 1990]

Cardinality: Possible Values [+/-].

A determiner D is cardinal iff $D \in \text{cardinal numbers} \geq 1$.

Genitive: Possible Values [+/-].

Possessive pronouns and the possessive morpheme (*'s*) are marked **gen**+; all other nouns are **gen**-.

Decreasing: Possible Values [+/-].

A set of Q properties is decreasing iff whenever $s \leq t$ and $t \in Q$ then $s \in Q$. A function f is decreasing iff for all properties $f(s)$ is a decreasing set.

A non-trivial NP (one with a Det) is decreasing iff its denotation in any model is decreasing. [Keenan and Stavi, 1986]

Constancy: Possible Values [+/-].

If A and B are sets denoting an NP and associated predicate, respectively, and E is a domain, then we say that a determiner displays constancy if $(A \cup B) \subseteq E \subseteq E'$ then $\text{Det}_EAB \leftrightarrow \text{Det}_{E'}AB$. Modified from [Partee *et al.*, 1990]

Complement: Possible Values [+/-].

A determiner Q is positive complement if and only if for every set X , there exists a continuous set of possible values for the size of the negated determined set, $\text{NOT}(QX)$, and the cardinality of QX is the only aspect of QX that can be negated. (adapted from [Mateyak, 1997])

The **wh**-feature has been discussed in the linguistics literature mainly in relation to wh-movement and with respect to NPs and nouns as well as determiners. We give a shallow but useful working definition of the **wh**-feature below:

Wh: Possible Values [+/-].

Interrogative determiners are **wh**+; all other determiners are **wh**-.

The **agr** feature is inherently a noun feature. While determiners are not morphologically marked for agreement in English many of them are sensitive to number. Many determiners are semantically either singular or plural and must adjoin to nouns which are the same. For example, *a* can only adjoin to singular nouns (*a dog* vs **a dogs* while *many* must have plurals (*many dogs* vs **many dog*). Other determiners such as *some* are unspecified for agreement in our analysis because they are compatible with either singulars or plurals (*some dog*, *some dogs*). The possible values of agreement for determiners are: [3sg, 3pl, 3].

Det	definite	quan	card	gen	wh	decreas	const	agr	compl
all	−	+	−	−	−	−	+	3pl	+
both	+	−	−	−	−	−	+	3pl	+
this	+	−	−	−	−	−	+	3sg	−
these	+	−	−	−	−	−	+	3pl	−
that	+	−	−	−	−	−	+	3sg	−
those	+	−	−	−	−	−	+	3pl	−
what	−	−	−	−	+	−	+	3	−
whatever	−	−	−	−	−	−	+	3	−
which	−	−	−	−	+	−	+	3	−
whichever	−	−	−	−	−	−	+	3	−
the	+	−	−	−	−	−	+	3	−
each	−	+	−	−	−	−	+	3sg	−
every	−	+	−	−	−	−	+	3sg	+
a/an	−	+	−	−	−	−	+	3sg	+
some ₁	−	+	−	−	−	−	+	3	−
some ₂	−	+	−	−	−	−	−	3pl	−
any	−	+	−	−	−	−	+	3sg	+
another	−	+	−	−	−	−	+	3sg	+
few	−	+	−	−	−	+	−	3pl	−
a few	−	+	−	−	−	−	+	3pl	−
many	−	+	−	−	−	−	−	3pl	+
many a/an	−	+	−	−	−	−	−	3sg	+
several	−	+	−	−	−	−	+	3pl	−
various	−	−	−	−	−	−	+	3pl	−
sundry	−	−	−	−	−	−	+	3pl	−
no	−	+	−	−	−	+	+	3	−
neither	−	−	−	−	−	+	+	3	−
either	−	−	−	−	−	−	+	3	−
GENITIVE	+	−	−	+	−	−	+	UN ²	−
CARDINAL	−	+	+	−	−	−	+	3pl ³	− ⁴
PARTITIVE	−	+/- ⁵	−	−	−	−	+	UN	+/-

Table 18.1: Determiner Features associated with D anchors

The determiner tree in Figure 18.2 shows the appropriate feature values for the determiner *these*, while Table 18.1 shows the corresponding feature values of several other common determiners.

In addition to the features that represent their own properties, determiners also have features to represent the selectional restrictions they impose on the NPs they take as complements. The selectional restriction features of a determiner appear on the NP footnode of the auxiliary tree that the determiner anchors. The NP_f node in Figure 18.2 shows the selectional feature restriction imposed by *these*⁶, while Table 18.2 shows the corresponding selectional feature restrictions of several other determiners.

18.1 The Wh-Feature

A determiner with a **wh+** feature is always the left-most determiner in linear order since no determiners have selectional restrictions that allow them to adjoin onto an NP with a +wh feature value. The presence of a wh+ determiner makes the entire NP wh+, and this is correctly represented by the coindexation of the determiner and root NP nodes' values for the wh-feature. Wh+ determiners' selectional restrictions on the NP foot node of their tree only allows them to adjoin onto NPs that are **wh-** or unspecified for the wh-feature. Therefore ungrammatical sequences such as **which what dog* are impossible. The adjunction of **wh+** determiners onto **wh+** pronouns is also prevented by the same mechanism.

18.2 Multi-word Determiners

The system recognizes the multi-word determiners *a few* and *many a*. The features for a multi-word determiner are located on the parent node of its two components (see Figure 18.3). We chose to represent these determiners as multi-word constituents because neither determiner retains the same set of features as either of its parts. For example, the determiner *a* is 3sg and *few* is decreasing, while *a few* is 3pl and increasing. Additionally, *many* is 3pl and *a* displays constancy, but *many a* is 3sg and does not display constancy. Example sentences appear in (270)-(271).

- Multi-word Determiners

(270) **a few** teaspoons of sugar should be adequate .

(271) **many a** man has attempted that stunt, but none have succeeded .

²We use the symbol UN to represent the fact that the selectional restrictions for a given feature are unspecified, meaning the noun phrase that the determiner selects can be either positive or negative for this feature.

³Except *one* which is 3sg.

⁴Except *one* which is **compl+**.

⁵A partitive can be either **quan+** or **quan-**, depending upon the nature of the noun that anchors the partitive. If the anchor noun is modified, then the quantity feature is determined by the modifier's quantity value.

⁶In addition to this tree, *these* would also anchor another auxiliary tree that adjoins onto **card+** determiners.

⁷*one* differs from the rest of CARD in selecting singular nouns

Det	defin	quan	card	gen	wh	decreas	const	agr	compl	e.g.
all	–	–	–	–	–	–	–	3pl	–	<i>dogs</i>
	+	–	–	UN	–	UN	UN	3pl	–	<i>these dogs</i>
	UN	UN	+	UN	UN	UN	UN	3pl	UN	<i>five dogs</i>
both	–	–	–	–	–	–	–	3pl	–	<i>dogs</i>
	+	–	–	UN	–	UN	UN	3pl	–	<i>these dogs</i>
this/that	–	–	–	–	–	–	–	3sg	–	<i>dog</i>
	–	+	UN	UN	–	+	–	3	UN	<i>few dogs</i>
	–	+	UN	UN	–	–	–	3pl	+	<i>many dogs</i>
	UN	UN	+	UN	UN	UN	UN	3sg	UN	<i>five dogs</i>
these/those	–	–	–	–	–	–	–	3pl	–	<i>dogs</i>
	–	+	UN	UN	–	+	–	3pl	UN	<i>few dogs</i>
	UN	UN	+	UN	UN	UN	UN	3pl	UN	<i>five dogs</i>
what/which whatever whichever	–	–	–	–	–	–	–	3	–	<i>dog(s)</i>
	–	+	UN	UN	–	+	–	3	UN	<i>few dogs</i>
	UN	UN	+	UN	UN	UN	UN	3	UN	<i>many dogs</i>
the	–	–	–	–	–	–	–	3	–	<i>dog(s)</i>
	–	+	UN	UN	–	+	–	3	UN	<i>few dogs</i>
	+	–	–	–	–	–	–	UN	–	<i>the me</i>
	–	+	UN	UN	–	–	–	3pl	+	<i>many dogs</i>
	UN	UN	+	UN	UN	UN	UN	3	UN	<i>five dogs</i>
every/each	–	–	–	–	–	–	–	3sg	–	<i>dog</i>
	–	+	UN	UN	–	+	–	3	UN	<i>few dogs</i>
	UN	UN	+	UN	UN	UN	UN	3	UN	<i>five dogs</i>
a/an	–	–	–	–	–	–	–	3sg	–	<i>dog</i>
some _{1,2} some ₁	–	–	–	–	–	–	–	3	–	<i>dog(s)</i>
	UN	UN	+	UN	UN	UN	UN	3pl	UN	<i>dogs</i>
any	–	–	–	–	–	–	–	3sg	–	<i>dog</i>
	–	+	UN	UN	–	+	–	3	UN	<i>few dogs</i>
	UN	UN	+	UN	UN	UN	UN	3	UN	<i>five dogs</i>
another	–	–	–	–	–	–	–	3sg	–	<i>dog</i>
	–	+	UN	UN	–	+	–	3	UN	<i>few dogs</i>
	UN	UN	+	UN	UN	UN	UN	3	UN	<i>five dogs</i>
few	–	–	–	–	–	–	–	3pl	–	<i>dogs</i>
a few	–	–	–	–	–	–	–	3pl	–	<i>dogs</i>
many	–	–	–	–	–	–	–	3pl	–	<i>dogs</i>
many a/an	–	–	–	–	–	–	–	3sg	–	<i>dog</i>
several	–	–	–	–	–	–	–	3pl	–	<i>dogs</i>
various	–	–	–	–	–	–	–	3pl	–	<i>dogs</i>
sundry	–	–	–	–	–	–	–	3pl	–	<i>dogs</i>
no	–	–	–	–	–	–	–	3	–	<i>dog(s)</i>
neither	–	–	–	–	–	–	–	3sg	–	<i>dog</i>
either	–	–	–	–	–	–	–	3sg	–	<i>dog</i>

Table 18.2: Selectional Restrictions Imposed by Determiners on the NP foot node

Det	definite	quan	card	gen	wh	decreas	const	agr	compl
GENITIVE	—	—	—	—	—	—	—	3	—
	—	+	UN	UN	—	+	—	3	UN
	—	+	UN	UN	—	—	—	3pl	+
	UN	UN	+	UN	UN	UN	UN	3	UN
	—	+	—	—	—	—	+	3pl	—
	—	—	—	—	—	—	+	3pl	—
CARDINAL	—	—	—	—	—	—	—	3pl ⁷	—
PARTITIVE	UN	UN	UN	UN	—	UN	UN	UN	UN

Table 18.3: Selectional Restrictions Imposed by Groups of Determiners/Determiner Constructions

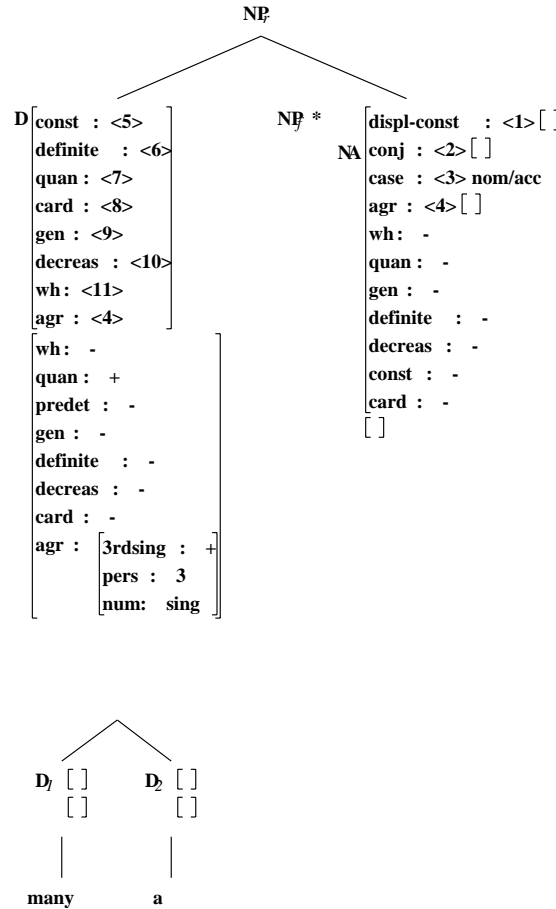


Figure 18.3: Multi-word Determiner tree: β DDnx

18.3 Genitive Constructions

There are two kinds of genitive constructions: genitive pronouns, and genitive NP's (which have an explicit genitive marker, 's, associated with them). It is clear from examples such as

her dog returned home and *her five dogs returned home* vs **dog returned home* that genitive pronouns function as determiners and as such, they sequence with the rest of the determiners. The features for the genitives are the same as for other determiners. Genitives are not required to agree with either the determiners or the nouns in the NPs that they modify. The value of the **agr** feature for an NP with a genitive determiner depends on the NP to which the genitive determiner adjoins. While it might seem to make sense to take *their* as 3pl, *my* as 1sg, and *Alfonso's* as 3sg, this number and person information only effects the genitive NP itself and bears no relationship to the number and person of the NPs with these items as determiners. Consequently, we have represented **agr** as unspecified for genitives in Table 18.1.

Genitive NP's are particularly interesting because they are potentially recursive structures. Complex NP's can easily be embedded within a determiner.

(272) [[[John]'s friend from high school]'s uncle]'s mother came to town.

There are two things to note in the above example. One is that in embedded NPs, the genitive morpheme comes at the end of the NP phrase, even if the head of the NP is at the beginning of the phrase. The other is that the determiner of an embedded NP can also be a genitive NP, hence the possibility of recursive structures.

In the XTAG grammar, the genitive marker, 's, is separated from the lexical item that it is attached to and given its own category (G). In this way, we can allow the full complexity of NP's to come from the existing NP system, including any recursive structures. As with the simple determiners, there is one auxiliary tree structure for genitives which adjoins to NPs. As can be seen in 18.4, this tree is anchored by the genitive marker 's and has a branching D node which accomodates the additional internal structure of genitive determiners. Also, like simple determiners, there is one initial tree structure (Figure 18.5) available for substitution where needed, as in, for example, the Determiner Gerund NP tree (see Chapter 17 for discussion on determiners for gerund NP's).

Since the NP node which is sister to the G node could also have a genitive determiner in it, the type of genitive recursion shown in (272) is quite naturally accounted for by the genitive tree structure used in our analysis.

18.4 Partitive Constructions

The deciding factor for including an analysis of partitive constructions(e.g. *some kind of, all of*) as a complex determiner constructions was the behavior of the agreement features. If partitive constructions are analyzed as an NP with an adjoined PP, then we would expect to get agreement with the head of the NP (as in (273)). If, on the other hand, we analyze them as a determiner construction, then we would expect to get agreement with the noun that the determiner sequence modifies (as we do in (274)).

(273) a *kind* [of these machines] *is* prone to failure.

(274) [a kind of] these *machines are* prone to failure.

In other words, for partitive constructions, the semantic head of the NP is the second rather than the first noun in linear order. That the agreement shown in (274) is possible suggests that

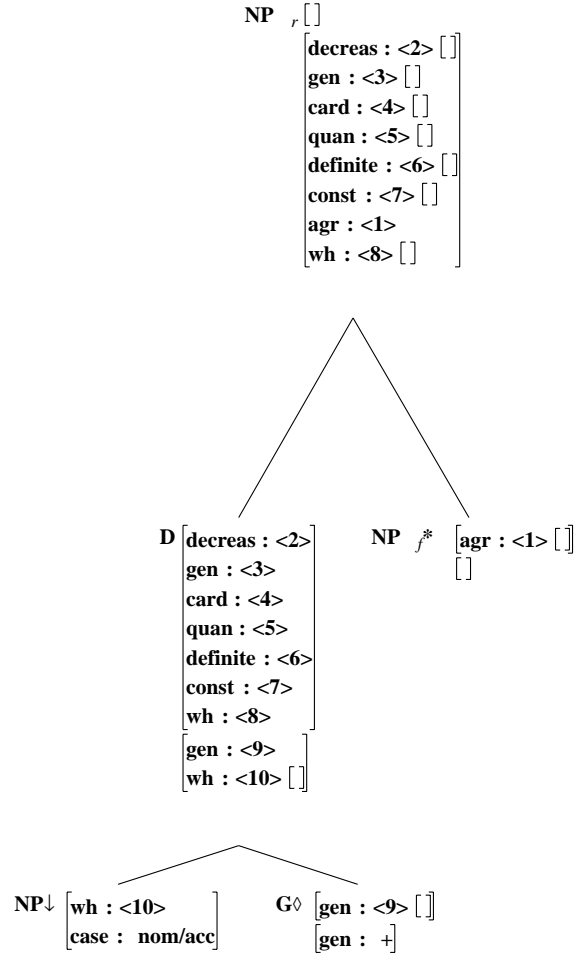


Figure 18.4: Genitive Determiner Tree

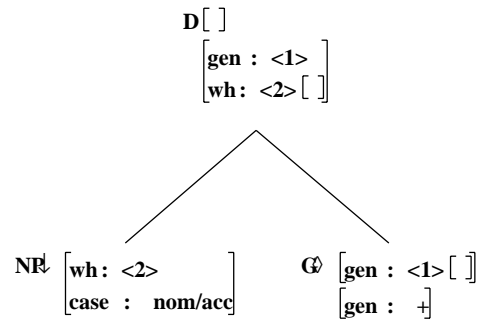


Figure 18.5: Genitive NP tree for substitution: α DnxG

the second noun in linear order in these constructions should also be treated as the syntactic head. Note that both the partitive and PP readings are usually possible for a particular NP. In

the cases where either the partitive or the PP reading is preferred, we take it to be just that, a preference, most appropriately modeled not in the grammar but in a component such as the heuristics used with the XTAG parser for reducing the analyses produced to the most likely.

In our analysis the partitive tree in Figure 18.6 is anchored by one of a limited group of nouns that can appear in the determiner portion of a partitive construction. A rough semantic characterization of these nouns is that they either represent quantity (e.g. *part, half, most, pot, cup, pound* etc.) or classification (e.g. *type, variety, kind, version* etc.). In the absence of a more implementable characterization we use a list of such nouns compiled from a descriptive grammar [Quirk *et al.*, 1985], a thesaurus, and from online corpora. In our grammar the nouns on the list are the only ones that select the partitive determiner tree.

Like other determiners, partitives can adjoin to an NP consisting of just a noun (*‘[a certain kind of] machine’*), or adjoin to NPs that already have determiners (*‘[some parts of] these machines’*). Notice that just as for the genitives, the complexity and the recursion are contained below the D node and rest of the structure is the same as for simple determiners.

18.5 Adverbs, Noun Phrases, and Determiners

Many adverbs interact with the noun phrase and determiner system in English. For example, consider sentences (275)-(282) below.

- (275) **Approximately** thirty people came to the lecture.
- (276) **Practically** every person in the theater was laughing hysterically during that scene.
- (277) **Only** John’s crazy mother can make stuffing that tastes so good.
- (278) **Relatively** few programmers remember how to program in COBOL.
- (279) **Not** every martian would postulate that all humans speak a universal language.
- (280) **Enough** money was gathered to pay off the group gift.
- (281) **Quite** a few burglaries occurred in that neighborhood last year.
- (282) I wanted to be paid **double** the amount they offered.

Although there is some debate in the literature as to whether these should be classified as determiners or adverbs, we believe that these items that interact with the NP and determiner system are in fact adverbs. These items exhibit a broader distribution than either determiners or adjectives in that they can modify many other phrasal categories, including adjectives, verb phrases, prepositional phrases, and other adverbs.

Using the determiner feature system, we can obtain a close approximation to an accurate characterization of the behavior of the adverbs that interact with noun phrases and determiners. Adverbs can adjoin to either a determiner or a noun phrase (see figure 18.7), with the adverbs restricting what types of NPs or determiners they can modify by imposing feature requirements on the foot D or NP node. For example, the adverb *approximately*, seen in (275) above, selects

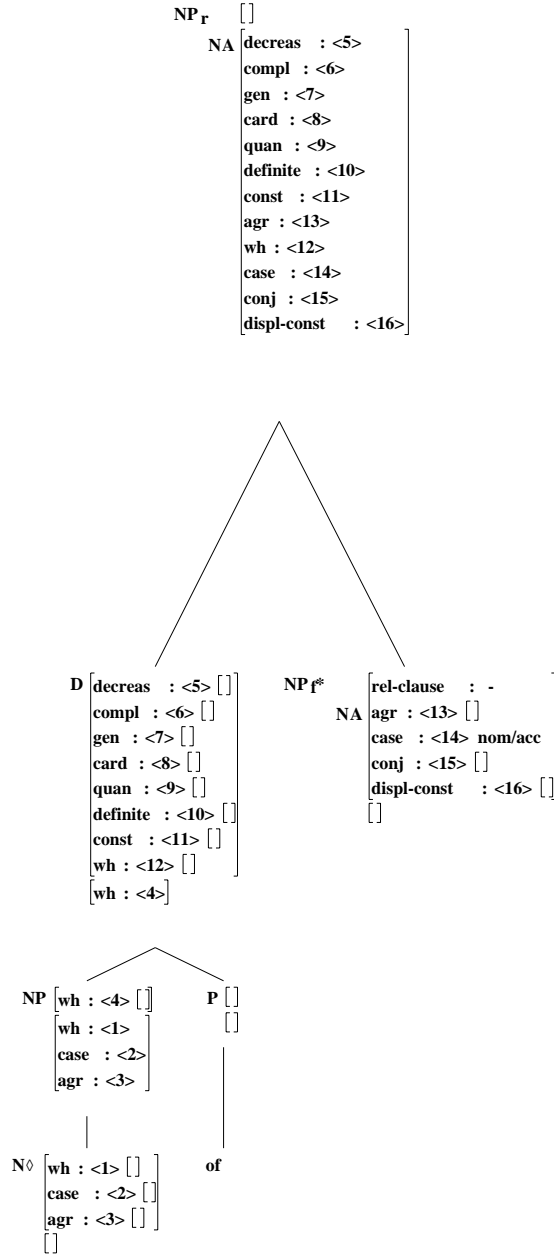


Figure 18.6: Partitive Determiner Tree

for determiners that are **card+**. The adverb *enough* in (280) is an example of an adverb that selects for a noun phrase, specifically a noun phrase that is not modified by a determiner.

Most of the adverbs that modify determiners and NPs divide into six classes, with some minor variation within classes, based on the pattern of these restrictions. Three of the classes are adverbs that modify determiners, while the other three modify NPs.

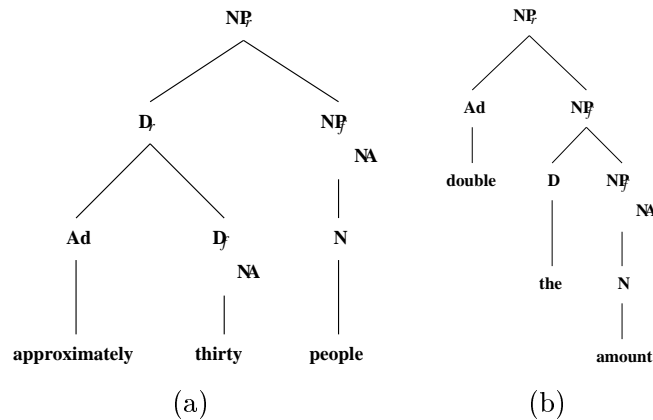


Figure 18.7: (a) Adverb modifying a determiner; (b) Adverb modifying a noun phrase

The largest of the five classes is the class of adverbs that modify cardinal determiners. This class includes, among others, the adverbs *about*, *at most*, *exactly*, *nearly*, and *only*. These adverbs have the single restriction that they must adjoin to determiners that are **card+**. Another class of adverbs consists of those that can modify the determiners *every*, *all*, *any*, and *no*. The adverbs in this class are *almost*, *nearly*, and *practically*. Closely related to this class are the adverbs *mostly* and *roughly*, which are restricted to modifying *every* and *all*, and *hardly*, which can only modify *any*. To select for *every*, *all*, and *any*, these adverbs select for determiners that are [**quan+**, **card-**, **const+**, **compl+**], and to select for *no*, the adverbs choose a determiner that is [**quan+**, **decreas+**, **const+**]. The third class of adverbs that modify determiners are those that modify the determiners *few* and *many*, representable by the feature sequences [**quan+**, **decreas+**, **const-**] and [**quan+**, **decreas-**, **const-**, **3pl**, **compl+**], respectively. Examples of these adverbs are *awfully*, *fairly*, *relatively*, and *very*.

Of the three classes of adverbs that modify noun phrases, one actually consists of a single adverb *not*, that only modifies determiners that are **compl+**. Another class consists of the focus adverbs, *at least*, *even*, *only*, and *just*. These adverbs select NPs that are **wh-** and **card-**. For the NPs that are **card+**, the focus adverbs actually modify the cardinal determiner, and so these adverbs are also included in the first class of adverbs mentioned in the previous paragraph. The last major class that modify NPs consist of the adverbs *double* and *twice*, which select NPs that are [**definite+**] (i.e., *the*, *this/that/those/these*, and the genitives).

Although these restrictions succeed in recognizing the correct determiner/adverb sequences, a few unacceptable sequences slip through. For example, in handling the second class of adverbs mentioned above, *every*, *all*, and *any* share the features [**quan+**, **card-**, **const+**, **compl+**] with *a* and *another*, and so **nearly a man* is acceptable in this system. In addition to this over-generation within a major class, the adverb *quite* selects for determiners and NPs in what seems to be a purely idiosyncratic fashion. Consider the following examples.

- (283) a. **Quite** a few members of the audience had to leave.
 b. There were **quite** many new participants at this year's conference.
 c. **Quite** few triple jumpers have jumped that far.

- d. Taking the day off was **quite** the right thing to do.
- e. The recent negotiation fiasco is **quite** another issue.
- f. Pandora is **quite** a cat!

In examples (283a)-(283c), *quite* modifies the determiner, while in (283d)-(283f), *quite* modifies the entire noun phrase. Clearly, it functions in a different manner in the two sets of sentences; in (283a)-(283c), *quite* intensifies the amount implied by the determiner, whereas in (283d)-(283f), it singles out an individual from the larger set to which it belongs. To capture the selectional restrictions needed for (283a)-(283c), we utilize the two sets of features mentioned previously for selecting *few* and *many*. However, *a few* cannot be singled out so easily; using the sequence [**quan+**, **card-**, **decreas-**, **const+**, **3pl**, **compl-**], we also accept the ungrammatical NPs **quite several members* and **quite some members* (where *quite* modifies *some*). In selecting *the* as in (d) with the features [**definite+**, **gen-**, **3sg**], *quite* also selects *this* and *that*, which are ungrammatical in this position. Examples (283e) and (283f) present yet another obstacle in that in selecting *another* and *a*, *quite* erroneously selects *every* and *any*.

It may be that there is an undiscovered semantic feature that would alleviate these difficulties. However, on the whole, the determiner feature system we have proposed can be used as a surprisingly efficient method of characterizing the interaction of adverbs with determiners and noun phrases.

Chapter 19

Modifiers

This chapter covers various types of modifiers: adverbs, prepositions, adjectives, and noun modifiers in noun-noun compounds.¹ These categories optionally modify other lexical items and phrases by adjoining onto them. In their modifier function these items are adjuncts; they are not part of the subcategorization frame of the items they modify. Examples of some of these modifiers are shown in (284)-(286).

(284) [_{ADV} certainly _{ADV}], the October 13 sell-off didn't settle any stomachs . (WSJ)

(285) Mr. Bakes [_{ADV} previously _{ADV}] had a turn at running Continental . (WSJ)

(286) most [_{ADJ} foreign _{ADJ}] [_N government _N] [_N bond _N] [prices] rose [_{PP} during the week _{PP}].

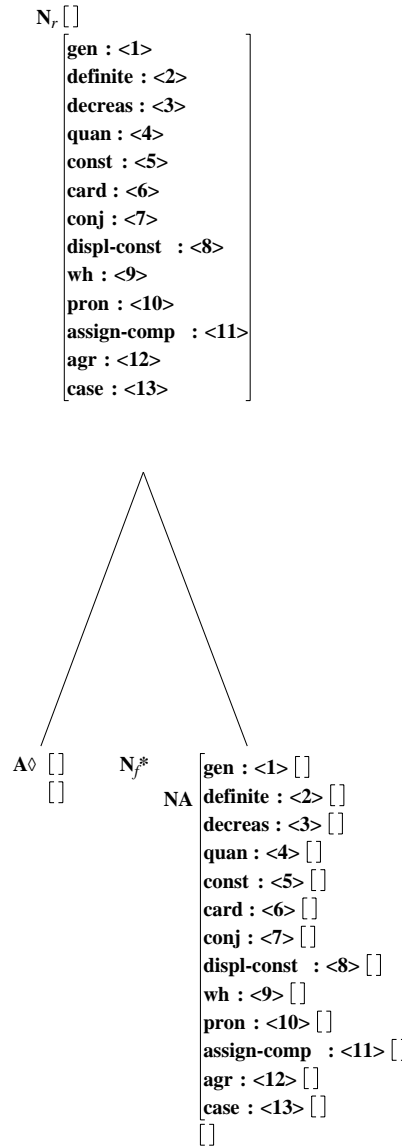
The trees used for the various modifiers are quite similar in form. The modifier anchors the tree and the root and foot nodes of the tree are of the category that the particular anchor modifies. Some modifiers, e.g. prepositions, select for their own arguments and those are also included in the tree. The foot node may be to the right or the left of the anchoring modifier (and its arguments) depending on whether that modifier occurs before or after the category it modifies. For example, almost all adjectives appear to the left of the nouns they modify, while prepositions appear to the right when modifying nouns.

19.1 Adjectives

In addition to being modifiers, adjectives in the XTAG English grammar can be also anchor clauses (see Adjective Small Clauses in Chapter 9). There is also one tree family, Intransitive with Adjective (Tnx0Vax1), that has an adjective as an argument and is used for sentences such as *Seth felt happy*. In that tree family the adjective substitutes into the tree rather than adjoining as is the case for modifiers.

As modifiers, adjectives anchor the tree shown in Figure 19.1. The features of the N node onto which the β An tree adjoins are passed through to the top node of the resulting N. The null adjunction marker (NA) on the N foot node imposes right binary branching such that each

¹Relative clauses are discussed in Chapter 14.

Figure 19.1: Standard Tree for Adjective modifying a Noun: βAn

subsequent adjective must adjoin on top of the leftmost adjective that has already adjoined. Due to the NA constraint, a sequence of adjectives will have only one derivation in the XTAG grammar. The adjective's morphological features such as superlative or comparative are instantiated by the morphological analyzer. See Chapter 22 for a description of how we handle comparatives. At this point, the treatment of adjectives in the XTAG English grammar does not include selectional or ordering restrictions. Consequently, any adjective can adjoin onto any noun and on top of any other adjective already modifying a noun. All of the modified noun

phrases shown in (287)-(290) currently parse with the same structure shown for *colorless green ideas* in Figure 19.2.

(287) big green bugs

(288) big green ideas

(289) colorless green ideas

(290) ?green big ideas

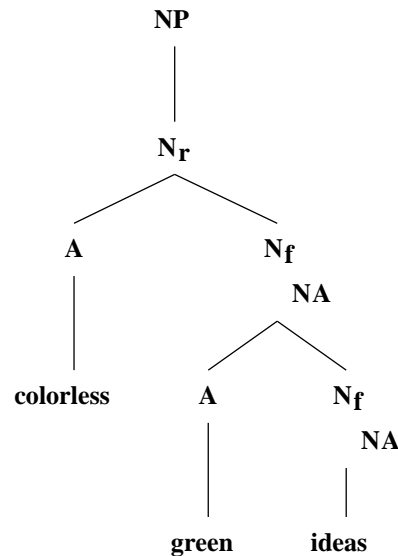


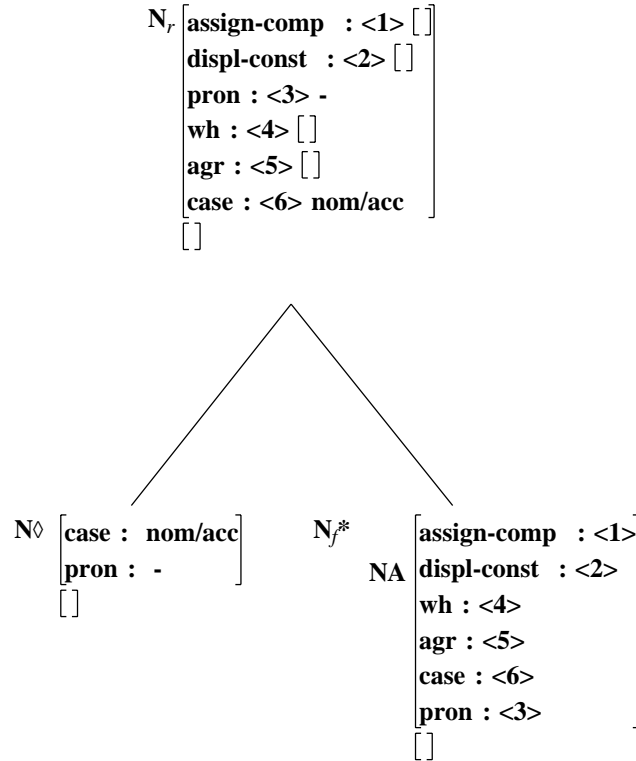
Figure 19.2: Multiple adjectives modifying a noun

While (288)-(290) are all semantically anomalous, (290) also suffers from an ordering problem that makes it seem ungrammatical in the absence of any licensing context. One could argue that the grammar should accept (287)-(289) but not (290). One of the future goals for the grammar is to develop a treatment of adjective ordering similar to that developed by [Hockey and Mateyak, 1998] for determiners². An adequate implementation of ordering restrictions for adjectives would rule out (290).

19.2 Noun-Noun Modifiers

Noun-noun compounding in the English XTAG grammar is very similar to adjective-noun modification. The noun modifier tree, shown in Figure 19.3, has essentially the same structure as the adjective modifier tree in Figure 19.1, except for the syntactic category label of the anchor.

²See Chapter 18 or [Hockey and Mateyak, 1998] for details of the determiner analysis.

Figure 19.3: Noun-noun compounding tree: βNn (not all features displayed)

Noun compounds have a variety of scope possibilities not available to adjectives, as illustrated by the single bracketing possibility in (291) and the two possibilities for (292). This ambiguity is manifested in the XTAG grammar by the two possible adjunction sites in the noun-noun compound tree itself. Subsequent modifying nouns can adjoin either onto the N_r node or onto the N anchor node of that tree, which results in exactly the two bracketing possibilities shown in (292). This inherent structural ambiguity results in noun-noun compounds regularly having multiple derivations. However, the multiple derivations are not a defect in the grammar because they are necessary to correctly represent the genuine ambiguity of these phrases.

(291) $[_N \text{ big } [_N \text{ green design } _N]_N]$

(292) $[_N \text{ computer } [_N \text{ furniture design } _N]_N]$
 $[_N [_N \text{ computer furniture } _N] \text{ design } _N]$

Noun-noun compounds have no restriction on number. XTAG allows nouns to be either singular or plural as in (293)-(295).

(293) Hyun is taking an algorithms course .

(294) waffles are in the frozen foods section .

(295) I enjoy the dog shows .

19.3 Time Noun Phrases

Although in general NPs cannot modify clauses or other NPs, there is a class of NPs, with meanings that relate to time, that can do so.³ We call this class of NPs “Time NPs”. Time NPs behave essentially like PPs. Like PPs, time NPs can adjoin at four places: to the right of an NP, to the right and left of a VP, and to the left of an S.

Time NPs may include determiners, as in *this month* in example (296), or may be single lexical items as in *today* in example (297). Like other NPs, time NPs can also include adjectives, as in example (301).

(296) Elvis left the building this week

(297) Elvis left the building today

(298) It has no bearing on our work force today (WSJ)

(299) The fire yesterday claimed two lives

(300) Today it has no bearing on our work force

(301) Michael late yesterday announced a buy-back program

The XTAG analysis for time NPs is fairly simple, requiring only the creation of proper NP auxiliary trees. Only nouns that can be part of time NPs will select the relevant auxiliary trees, and so only these type of NPs will behave like PPs under the XTAG analysis. Currently, about 60 words select Time NP trees, but since these words can form NPs that include determiners and adjectives, a large number of phrases are covered by this class of modifiers.

Corresponding to the four positions listed above, time NPs can select one of the four trees shown in Figure 19.4.

Determiners can be added to time NPs by adjunction in the same way that they are added to NPs in other positions. The trees in Figure 19.5 show that the structures of examples (296) and (297) differ only in the adjunction of *this* to the time NP in example (296).

The sentence

(302) Esso said the Whiting field started production Tuesday (WSJ)

has (at least) two different interpretations, depending on whether *Tuesday* attaches to *said* or to *started*. Valid time NP analyses are available for both these interpretations and are shown in Figure 19.6.

Derived tree structures for examples (298) – (301), which show the four possible time NP positions are shown in Figures 19.7 and 19.8. The derivation tree for example (301) is also shown in Figure 19.8.

³There may be other classes of NPs, such as directional phrases, such as *north*, *south* etc., which behave similarly. We have not yet analyzed these phrases.

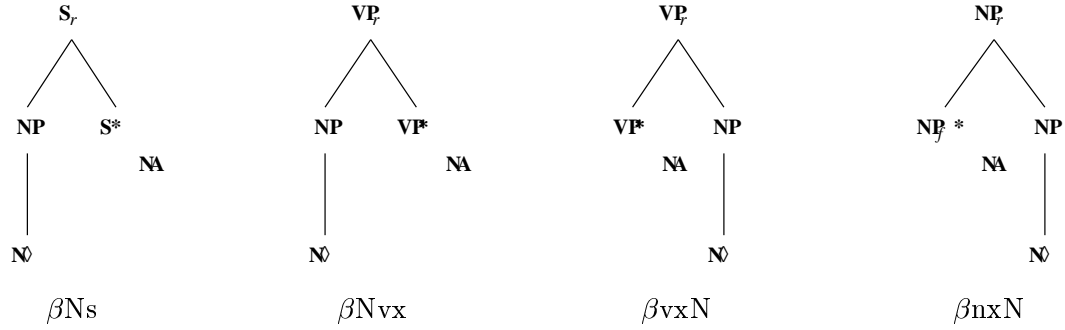
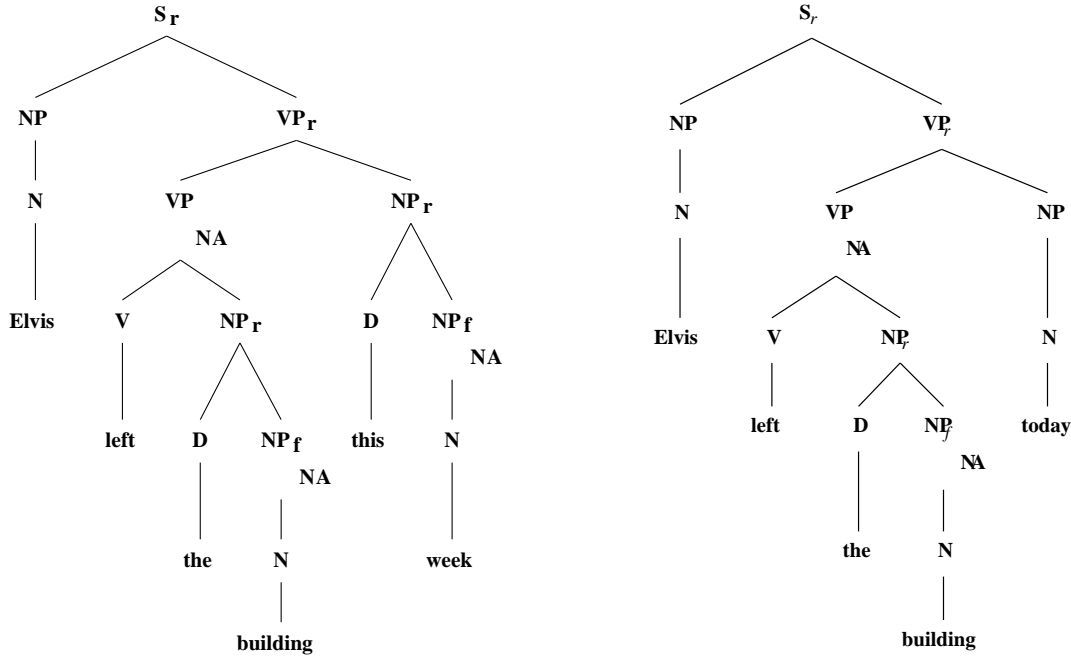

 Figure 19.4: Time Phrase Modifier trees: βNs , βNvx , βvxN , βNxN


Figure 19.5: Time NPs with and without a determiner

19.4 Prepositions

There are three basic types of prepositional phrases, and three places at which they can adjoin. The three types of prepositional phrases are: Preposition with NP Complement, Preposition with Sentential Complement, and Exhaustive Preposition. The three places are to the right of an NP, to the right of a VP, and to the left of an S. Each of the three types of PP can adjoin at each of these three places, for a total of nine PP modifier trees. Table 19.1 gives the tree family names for the various combinations of type and location. (See Section 23.4.2 for discussion of the $\beta spuPnx$, which handles post-sentential comma-separated PPs.)

The subset of preposition anchored modifier trees in Figure 19.9 illustrates the locations and the four PP types. Example sentences using the trees in Figure 19.9 are shown in (303)-(306). There are also more trees with multi-word prepositions as anchors. Examples of these are:

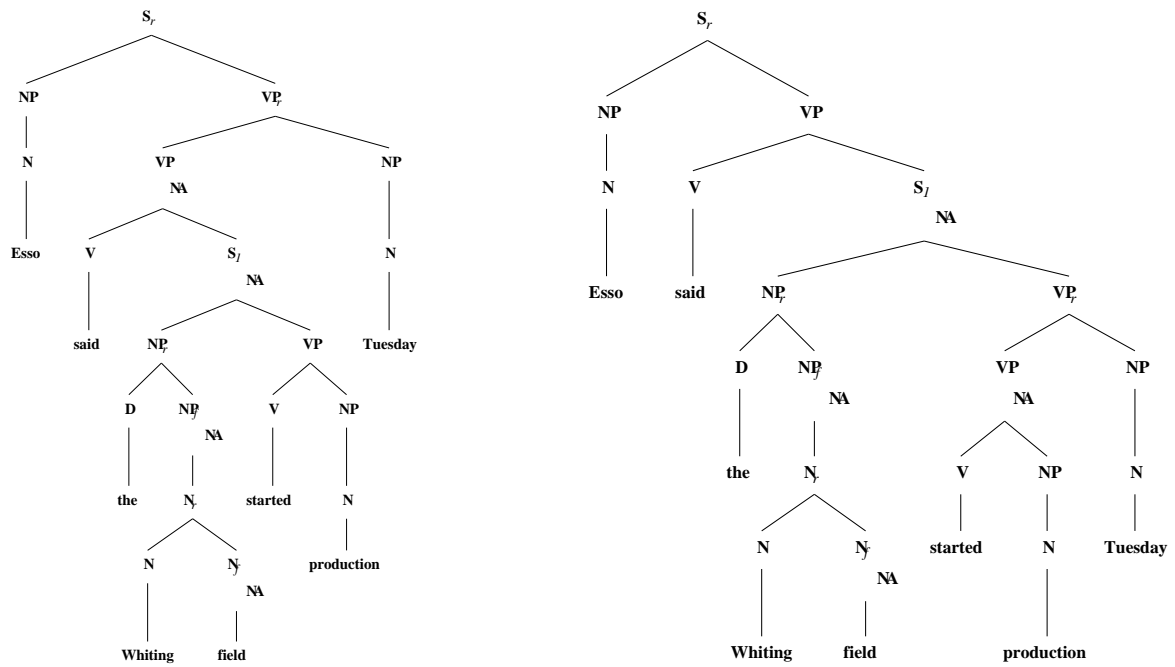


Figure 19.6: Time NP trees: Two different attachments

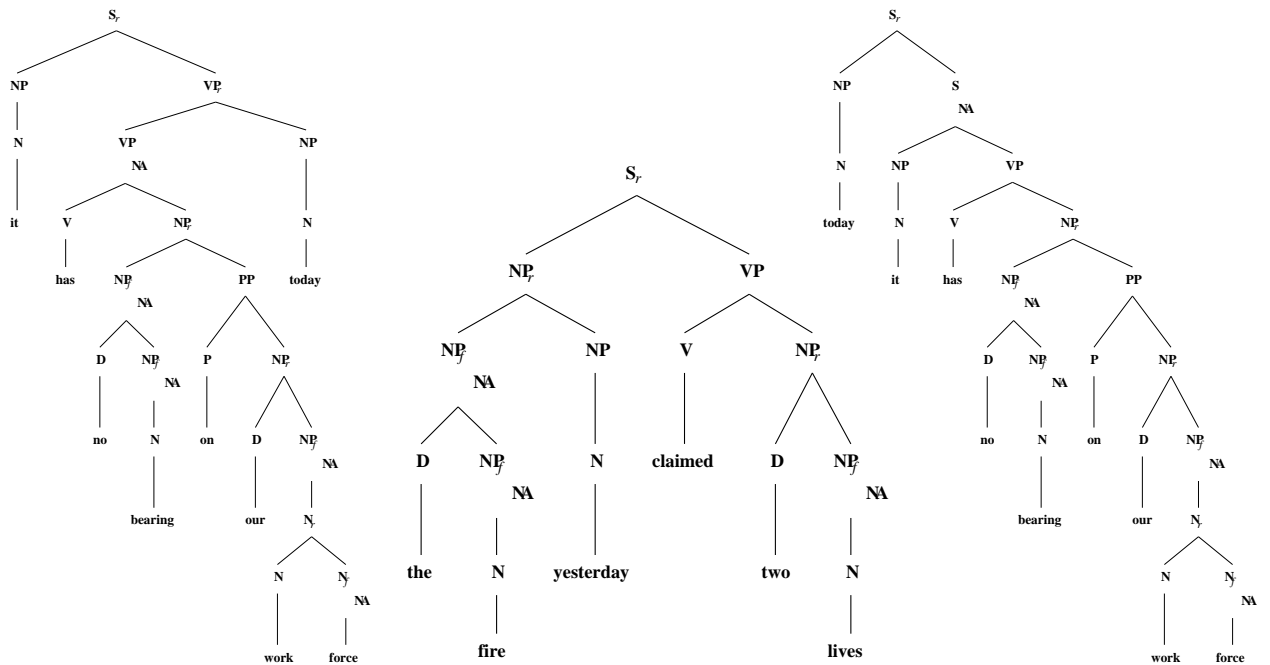
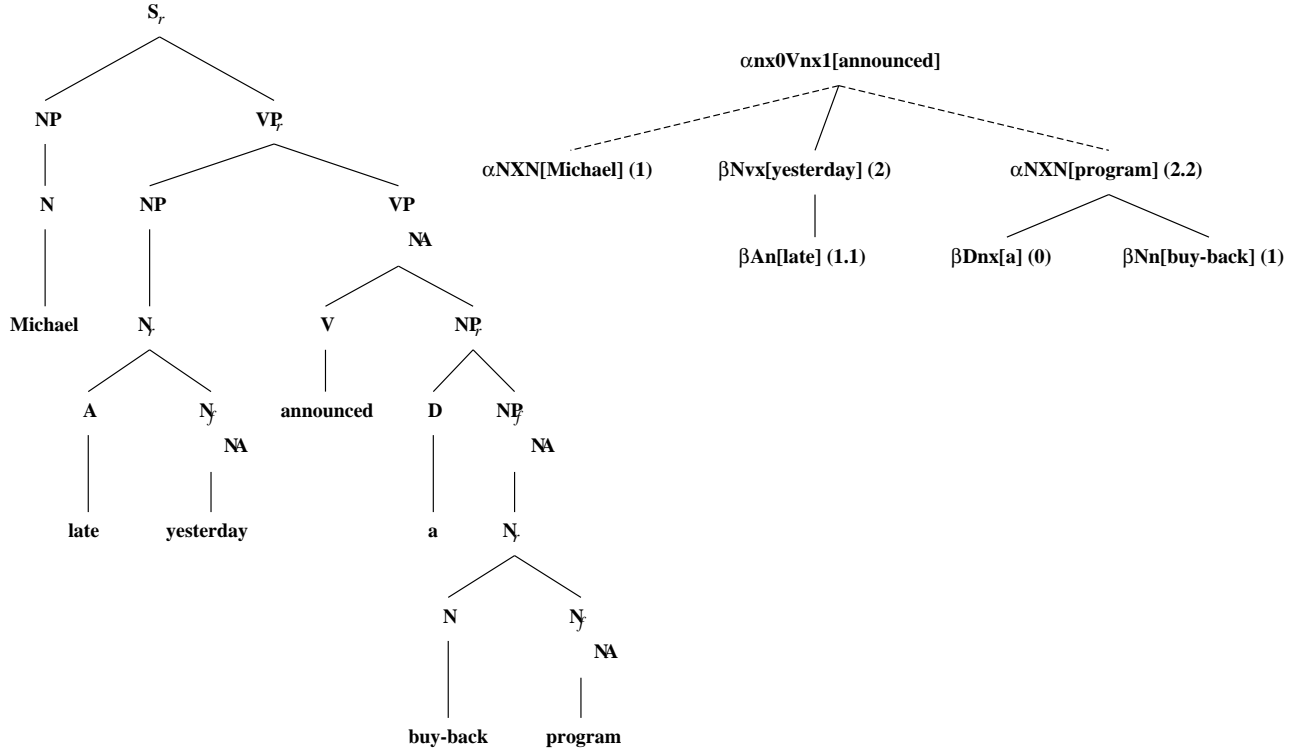


Figure 19.7: Time NPs in different positions (βvxN , βnxN and βNs)

ahead of, contrary to, at variance with and as recently as.


 Figure 19.8: Time NPs: Derived tree and Derivation (βNvx position)

	position and category modified		
	pre-sentential S modifier	post-NP NP modifier	post-VP VP modifier
S-complement	βPss	$\beta nxPs$	$\beta vxPs$
NP-complement	$\beta Pnxs$	$\beta nxPnx$	$\beta vxPnx$
no complement (exhaustive)	βPs	βnxP	βvxP

Table 19.1: Preposition Anchored Modifiers

- (303) [$_{PP}$ with Clove healthy $_{PP}$], the veterinarian's bill will be more affordable . (βPss^4)
- (304) The frisbee [$_{PP}$ in the brambles $_{PP}$] was hidden . ($\beta nxPnx$)
- (305) Clove played frisbee [$_{PP}$ outside $_{PP}$] . (βvxP)
- (306) Clove played frisbee [$_{PP}$ outside of the house $_{PP}$] . ($\beta vxPPnx$)

Prepositions that take NP complements assign accusative case to those complements (see section 4.4.3.1 for details). Most prepositions take NP complements. Subordinating conjunctions are analyzed in XTAG as Preps (see Section 15 for details). Additionally, a few non-conjunction prepositions take S complements (see Section 8.8 for details).

⁴ *Clove healthy* is an adjective small clause

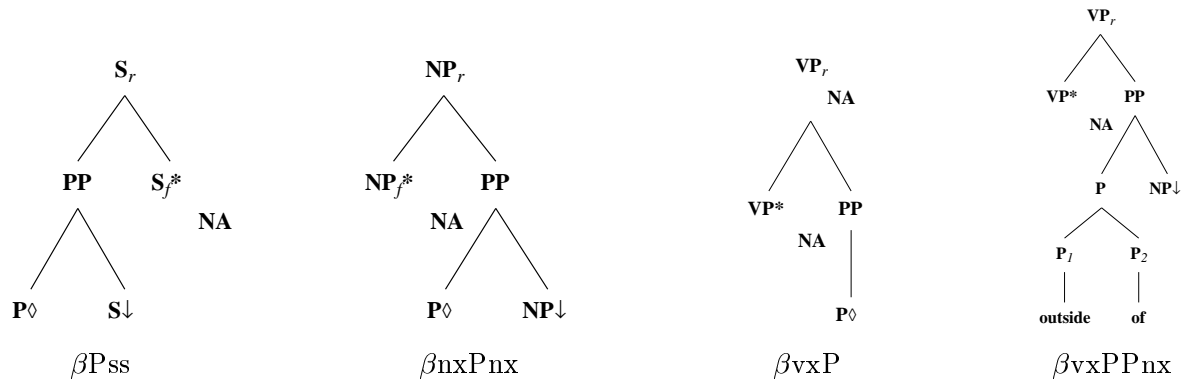


Figure 19.9: Selected Prepositional Phrase Modifier trees: βP_{ss} , $\beta n x P_{n x}$, $\beta v x P$ and $\beta v x P P_{n x}$

19.5 Adverbs

In the English XTAG grammar, VP and S-modifying adverbs anchor the auxiliary trees $\beta A R B$ s, $\beta s A R B$, $\beta v x A R B$ and $\beta A R B v x$,⁵ allowing pre and post modification of S's and VP's. Besides the VP and S-modifying adverbs, the grammar includes adverbs that modify other categories. Examples of adverbs modifying an adjective, an adverb, a PP, an NP, and a determiner are shown in (307)-(314). (See Sections 23.1.5 and 23.4.1 for discussion of the $\beta p u A R B p u v x$ and $\beta s p u A R B$, which handle pre-verbal parenthetical adverbs and post-sentential comma-separated adverbs.)

- Modifying an adjective

(307) **extremely** good

(308) **rather** tall

(309) rich **enough**

- Modifying an adverb

(310) oddly **enough**

(311) **very** well

- Modifying a PP

(312) **right** through the wall

- Modifying a NP

(313) **quite** some time

⁵In the naming conventions for the XTAG trees, ARB is used for adverbs. Because the letters in A, Ad, and Adv are all used for other parts of speech (adjective, determiner and verb), ARB was chosen to eliminate ambiguity. Appendix D contains a full explanation of naming conventions.

- Modifying a determiner

(314) **exactly** five men

XTAG has separate trees for each of the modified categories and for pre and post modification where needed. The kind of treatment given to adverbs here is very much in line with the base-generation approach proposed by [Ernst, 1983], which assumes all positions where an adverb can occur to be base-generated, and that the semantics of the adverb specifies a range of possible positions occupied by each adverb. While the relevant semantic features of the adverbs are not currently implemented, implementation of semantic features is scheduled for future work. The trees for adverb anchored modifiers are very similar in form to the adjective anchored modifier trees. Examples of two of the basic adverb modifier trees are shown in Figure 19.10.

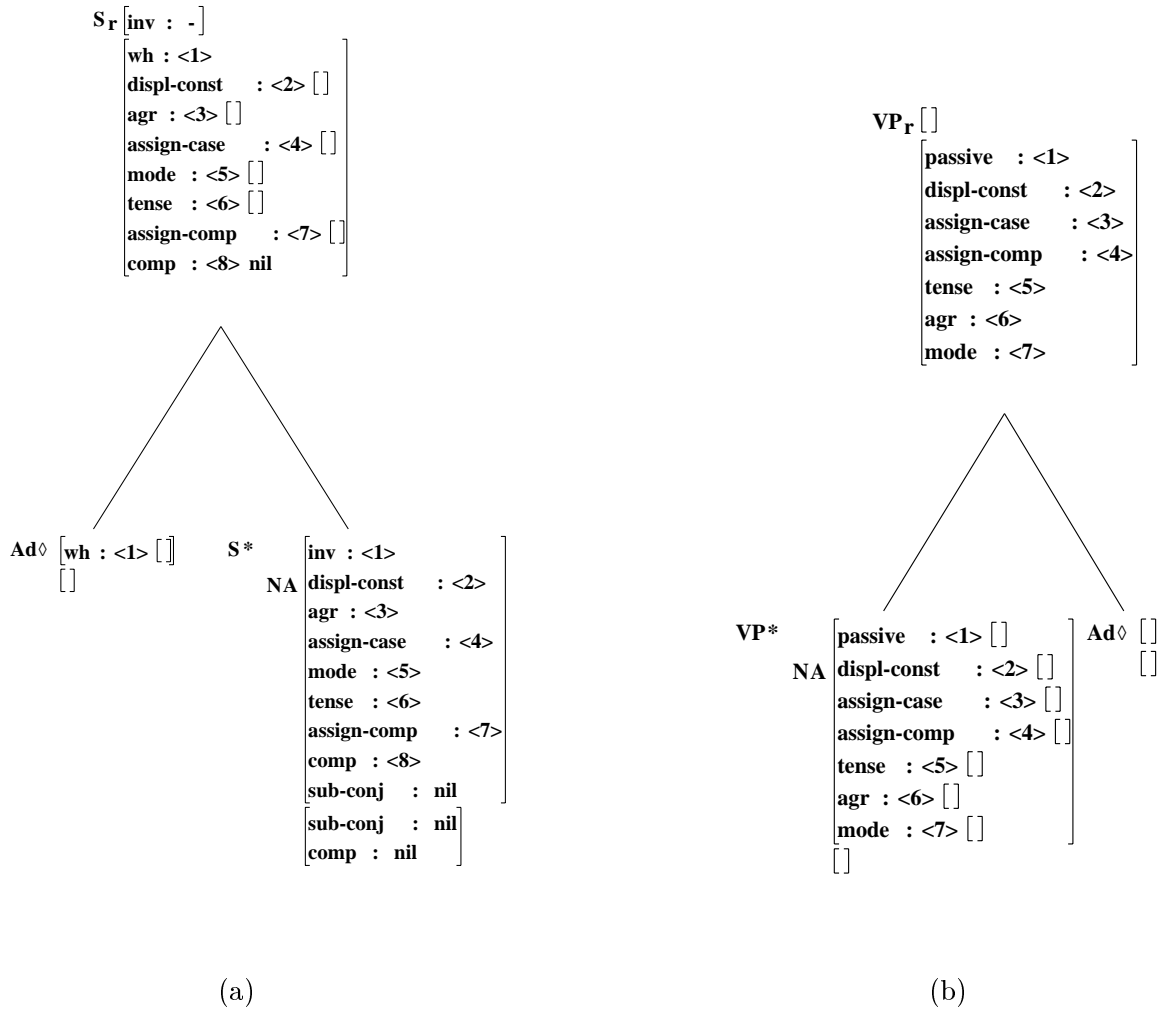


Figure 19.10: Adverb Trees for pre-modification of S: β ARBs (a) and post-modification of a VP: β_{vx} ARB (b)

Like the adjective anchored trees, these trees also have the NA constraint on the foot node to restrict the number of derivations produced for a sequence of adverbs. Features of the modified category are passed from the foot node to the root node, reflecting correctly that these types of properties are unaffected by the adjunction of an adverb. A summary of the categories modified and the position of adverbs is given in Table 19.2.

Category Modified	Position with respect to item modified	
	Pre	Post
S	β ARBs	β sARB
VP	β ARBvx, β puARBpuvx	β vxARB
A	β ARBa	β aARB
PP	β ARBpx	β pxARB
ADV	β ARBarb	β arbARB
NP	β ARBnx	
Det	β ARBd	

Table 19.2: Simple Adverb Anchored Modifiers

In the English XTAG grammar, no traces are posited for wh-adverbs, in-line with the base-generation approach ([Ernst, 1983]) for various positions of adverbs. Since convincing arguments have been made against traces for adjuncts of other types (e.g. [Baltin, 1989]), and since the reasons for wanting traces do not seem to apply to adjuncts, we make the general assumption in our grammar that adjuncts do not leave traces. Sentence initial wh-adverbs select the same auxiliary tree used for other sentence initial adverbs (β ARBs) with the feature $\langle \mathbf{wh} \rangle = +$. Under this treatment, the derived tree for the sentence *How did you fall?* is as in Figure (19.11), with no trace for the adverb.

There is one more adverb modifier tree in the grammar which is not included in Table 19.2. This tree, shown in Figure 19.12, has a complex adverb phrase and is used for wh+ two-adverb phrases that occur sentence initially, such as in sentence (315). Since *how* is the only wh+ adverb, it is the only adverb that can anchor this tree.

(315) how quickly did Srini fix the problem ?

Focus adverbs such as *only*, *even*, *just* and *at least* are also handled by the system. Since the syntax allows focus adverbs to appear in practically any position, these adverbs select most of the trees listed in Table 19.2. It is left up to the semantics or pragmatics to decide the correct scope for the focus adverb for a given instance. In terms of the ability of the focus adverbs to modify at different levels of a noun phrase, the focus adverbs can modify either cardinal determiners or noun-cardinal noun phrases, and cannot modify at the level of noun. The tree for adverbial modification of noun phrases is in shown Figure 19.13(a).

In addition to *at least*, the system handles the other two-word adverbs, *at most* and *up to*, and the three-word *as-as* adverb constructions, where an adjective substitutes between the two occurrences of *as*. An example of a three-word *as-as* adverb is *as little as*. Except for the ability of *at least* to modify many different types of constituents as noted above, the multi-word adverbs are restricted to modifying cardinal determiners. Example sentences using the trees in Figure 19.13 are shown in (316)-(320).

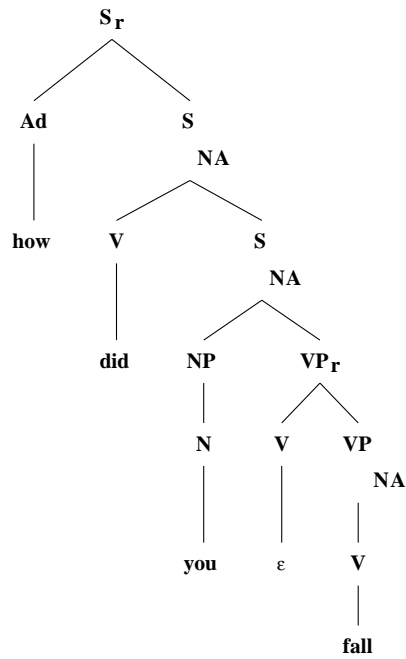


Figure 19.11: Derived tree for *How did you fall?*

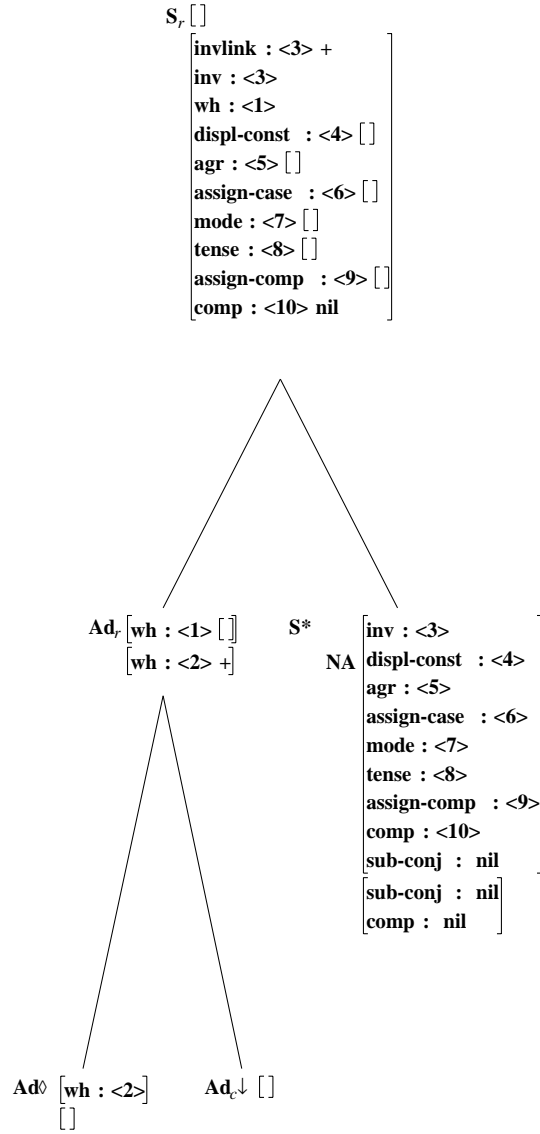


Figure 19.12: Complex adverb phrase modifier: β ARBarbs

- Focus Adverb modifying an NP

(316) **only** a member of our crazy family could pull off that kind of a stunt .

(317) **even** a flying saucer sighting would seem interesting in comparison with your story .

(318) The report includes a proposal for **at least** a partial impasse in negotiations .

- Multi-word adverbs modifying cardinal determiners

(319) **at most** ten people came to the party .

(320) They gave monetary gifts of **as little as** five dollars .

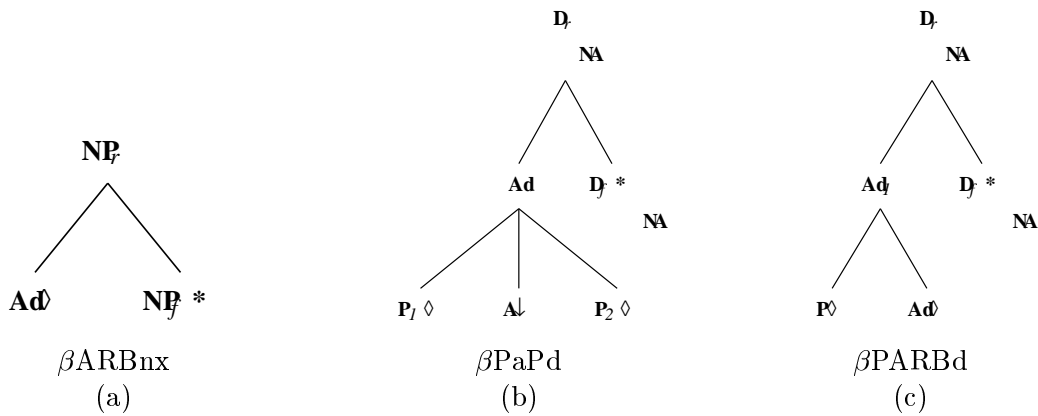


Figure 19.13: Selected Focus and Multi-word Adverb Modifier trees: $\beta\text{ARB}_{\text{nx}}$, βPARBd and βPaPd

The grammar also includes auxiliary trees anchored by multi-word adverbs like *a little*, *a bit*, *a mite*, *sort of*, *kind of*, etc..

Multi-word adverbs like *sort of* and *kind of* can pre- modify almost any non-clausal category. The only strict constraint on their occurrence is that they can't modify nouns (in which case an adjectival interpretation would obtain)⁶. The category which they scope over can be directly determined from their position, except for when they occur sentence finally in which case they are assumed to modify VP's. The complete list of auxiliary trees anchored by these adverbs are as follows: βNPax , βNPpx , βNPnx , βNPvx , βvxNP , βNParb . Selected trees are shown in Figure 19.14, and some examples are given in (321)-(324).

(321) John is **sort of** [_{AP} tired].

(322) John is **sort of** [_{PP} to the right].

⁶Note that there are semantic/lexical constraints even for the categories that these adverbs *can* modify, and no doubt invite a more in-depth analysis.

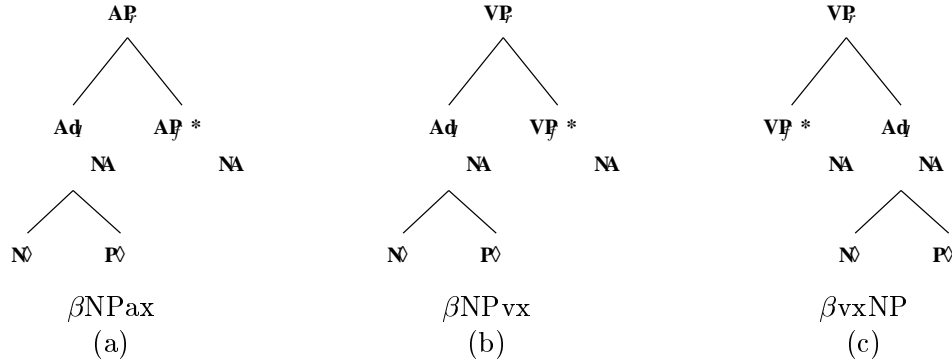


Figure 19.14: Selected Multi-word Adverb Modifier trees (for adverbs like *sort of*, *kind of*): βNPax , βNPvx , βvxNP .

(323) John could have been **sort of** [_{VP} eating the cake].

(324) John has been eating his cake **sort of** [_{ADV} slowly].

There are some multi-word adverbs that are, however, not so free in their distribution. Adverbs like *a little*, *a bit*, *a mite* modify AP's in predicative constructions (sentences with the copula and small clauses, AP complements in sentences with raising verbs, and AP's when they are subcategorized for by certain verbs (e.g., *John felt angry*). They can also post-modify VP's and PP's, though not as freely as AP's⁷. Finally, they also function as downtoners for almost all adverbials⁸. Some examples are provided in (325)-(328).

(325) Mickey is **a little** [_{AP} tired].

(326) The medicine [_{VP} has eased John's pain] **a little**.

(327) John is **a little** [_{PP} to the right].

(328) John has been reading his book **a little** [_{ADV} loudly].

Following their behavior as described above, the auxiliary trees they anchor are βDAax , βDApx , βvxDA , βDAarb , βDNax , βDNpx , βvxDN , βDNarb . Selected trees are shown in Figure 19.15).

19.6 Locative Adverbial Phrases

Locative adverbial phrases are multi-word adverbial modifiers whose meanings relate to spatial location. Locatives consist of a locative adverb (such as *ahead* or *downstream*) preceded by

⁷They can also appear before NP's, as in, "John wants *a little* sugar". However, here they function as multi-word determiners and should not be analyzed as adverbs.

⁸It is to be noted that this analysis, which allows these multiword adverbs to modify adjectival phrases as well as adverbials, will yield (not necessarily desirable) multiple derivations for a sentence like *John is a little unnecessarily stupid*. In one derivation, *a little* modifies the AP and in the other case, it modifies the adverb.

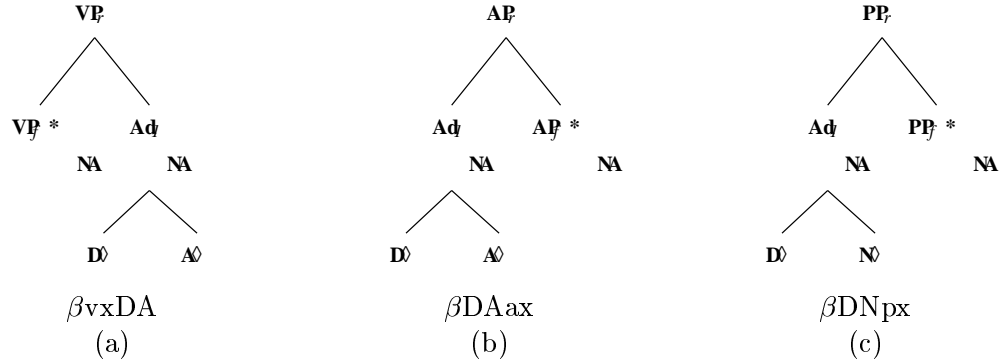


Figure 19.15: Selected Multi-word Adverb Modifier trees (for adverbs like *a little*, *a bit*): $\beta_{vx}DA$, $\beta_{DA}ax$, $\beta_{DN}px$.

an NP, an adverb, or nothing, as in Examples (329)–(331) respectively. The modifier as a whole describes a position relative to one previously specified in the discourse. The nature of the relation, which is usually a direction, is specified by the anchoring locative adverb (*behind*, *east*). If an NP or a second adverb is present in the phrase, it specifies the degree of the relation (for example: *three city blocks*, *many meters*, and *far*).

(329) The accident *three blocks ahead* stopped traffic

(330) The ship sank *far offshore*

(331) The trouble *ahead* distresses me

Locatives can modify NPs, VPs and Ss. They modify NPs only by right-adjoining post-positively, as in Example (329). Post-positive is also the more common position when a locative modifies either of the other categories. Locatives pre-modify VPs only when separated by balanced punctuation (commas or dashes). The trees locatives select when modifying NPs are shown in Figure 19.16.

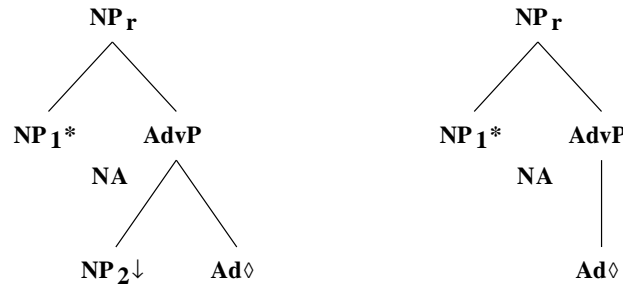


Figure 19.16: Locative Modifier Trees: $\beta_{nxxn}ARB$, $\beta_{nx}ARB$

When the locative phrase consists of only the anchoring locative adverb, as in Example (330), it uses the $\beta_{nx}ARB$ tree, shown in Figure 19.16, and its VP analogue, $\beta_{vx}ARB$. In addition, these are the trees selected when the locative anchor is modified by an adverb expressing degree,

as in Example 330. The degree adverb adjoins on to the anchor using the β ARBarb tree, which is described in Section 19.5. Figure 19.17 shows an example of these trees in action. Though there is a tree for a pre-sentential locative phrase, β nxARBs, there is no corresponding post-sentential tree, as it is highly debatable whether the post-sentential version actually has the entire sentence or just the preceding verb phrase as its scope. Thus, in accordance with XTAG practice, which considers ambiguous post-sentential modifiers to be VP-modifiers rather than S-modifiers, there is only a β vxnARB tree, as shown in Figure 19.17.

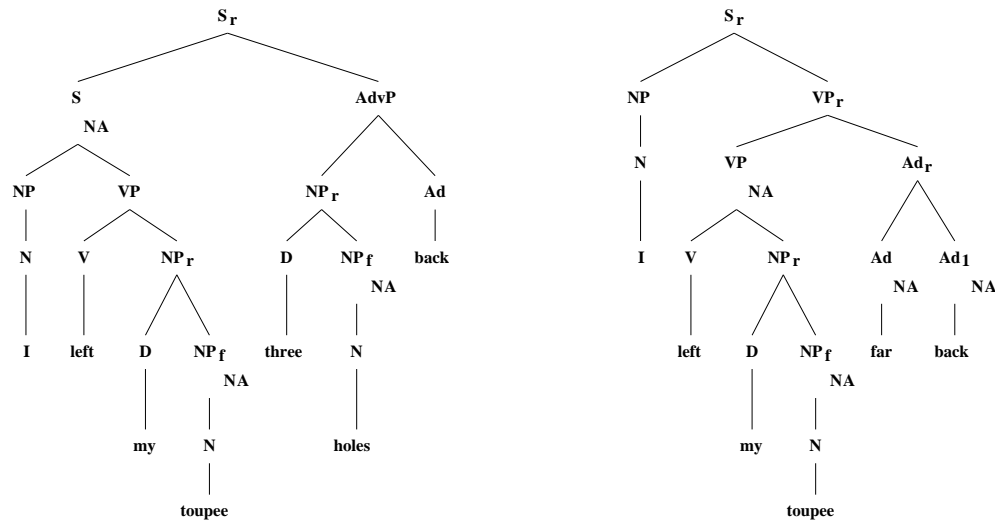


Figure 19.17: Locative Phrases featuring NP and Adverb Degree Specifications

One possible analysis of locative phrases with NPs might maintain that the NP is the head, with the locative adverb modifying the NP. This is initially attractive because of the similarity to time NPs, which also feature NPs that can modify clauses. This analysis seems insufficient, however, in light of the fact that virtually any NP can occur in locative phrases, as in example (332). Therefore, in the XTAG analysis the locative adverb anchors the locative phrase trees. A complete summary of all trees selected by locatives is contained in Table 19.3. 26⁹ adverbs select the locative trees.

(332) I left my toupee and putter *three holes back*

⁹Though nearly all of these adverbs are spatial in nature, this number also includes a few temporal adverbs, such as *ago*, that also select these trees.

Category Modified	Degree Phrase Type	
	NP	Ad/None
NP	β_{nxnxARB}	β_{nxARB}
VP (post)	β_{vxnxARB}	β_{vxARB}
VP (pre, punct-separated)	$\beta_{\text{punxARBpuvx}}$	$\beta_{\text{puARBpuvx}}$
S	β_{nxARBs}	β_{ARBs}

Table 19.3: Locative Modifiers

Chapter 20

Auxiliaries

Although there has been some debate about the lexical category of auxiliaries, the English XTAG grammar follows [McCawley, 1988], [Haegeman, 1991], and others in classifying auxiliaries as verbs. The category of verbs can therefore be divided into two sets, main or lexical verbs, and auxiliary verbs, which can co-occur in a verbal sequence. Only the highest verb in a verbal sequence is marked for tense and agreement regardless of whether it is a main or auxiliary verb. Some auxiliaries (*be*, *do*, and *have*) share with main verbs the property of having overt morphological marking for tense and agreement, while the modal auxiliaries do not. However, all auxiliary verbs differ from main verbs in several crucial ways.

- Multiple auxiliaries can occur in a single sentence, while a matrix sentence may have at most one main verb.
- Auxiliary verbs cannot occur as the sole verb in the sentence, but must be followed by a main verb.
- All auxiliaries precede the main verb in verbal sequences.
- Auxiliaries do not subcategorize for any arguments.
- Auxiliaries impose requirements on the morphological form of the verbs that immediately follow them.
- Only auxiliary verbs invert in questions (with the sole exception in American English of main verb *be*¹).
- An auxiliary verb must precede sentential negation (e.g. **John not goes*).
- Auxiliaries can form contractions with subjects and negation (e.g. *he'll*, *won't*).

The restrictions that an auxiliary verb imposes on the succeeding verb limits the sequence of verbs that can occur. In English, sequences of up to five verbs are allowed, as in sentence (333).

¹Some dialects, particularly British English, can also invert main verb *have* in yes/no questions (e.g. *have you any Grey Poupon* ?). This is usually attributed to the influence of auxiliary *have*, coupled with the historic fact that English once allowed this movement for all verbs.

(333) The music should have been being played [for the president] .

The required ordering of verb forms when all five verbs are present is:

modal base perfective progressive passive

The rightmost verb is the main verb of the sentence. While a main verb subcategorizes for the arguments that appear in the sentence, the auxiliary verbs select the particular morphological forms of the verb to follow each of them. The auxiliaries included in the English XTAG grammar are listed in Table 20.1 by type. The third column of Table 20.1 lists the verb forms that are required to follow each type of auxiliary verb.

TYPE	LEX ITEMS	SELECTS FOR
modals	<i>can, could, may, might, will, would, ought, shall, should need</i>	base form ² (e.g. <i>will go, might come</i>)
perfective	<i>have</i>	past participle (e.g. <i>has gone</i>)
progressive	<i>be</i>	gerund (e.g. <i>is going, was coming</i>)
passive	<i>be</i>	past participle (e.g. <i>was helped by Jane</i>)
do support	<i>do</i>	base form (e.g. <i>did go, does come</i>)
infinitive to	<i>to</i>	base form (e.g. <i>to go, to come</i>)

Table 20.1: Auxiliary Verb Properties

20.1 Non-inverted sentences

This section and the sections that follow describe how the English XTAG grammar accounts for properties of the auxiliary system described above.

In our grammar, auxiliary trees are added to the main verb tree by adjunction. Figure 20.1 shows the adjunction tree for non-inverted sentences.³

The restrictions outlined in column 3 of Table 20.1 are implemented through the features <**mode**>, <**perfect**>, <**progressive**> and <**passive**>. The syntactic lexicon entries for the auxiliaries gives values for these features on the foot node (VP*) in Figure 20.1. Since the top features of the foot node must eventually unify with the bottom features of the node it adjoins

²There are American dialects, particularly in the South, which allow double modals such as *might could* and *might should*. These constructions are not allowed in the XTAG English grammar.

³We saw this tree briefly in section 4.4.3.2, but with most of its features missing. The full tree is presented here.

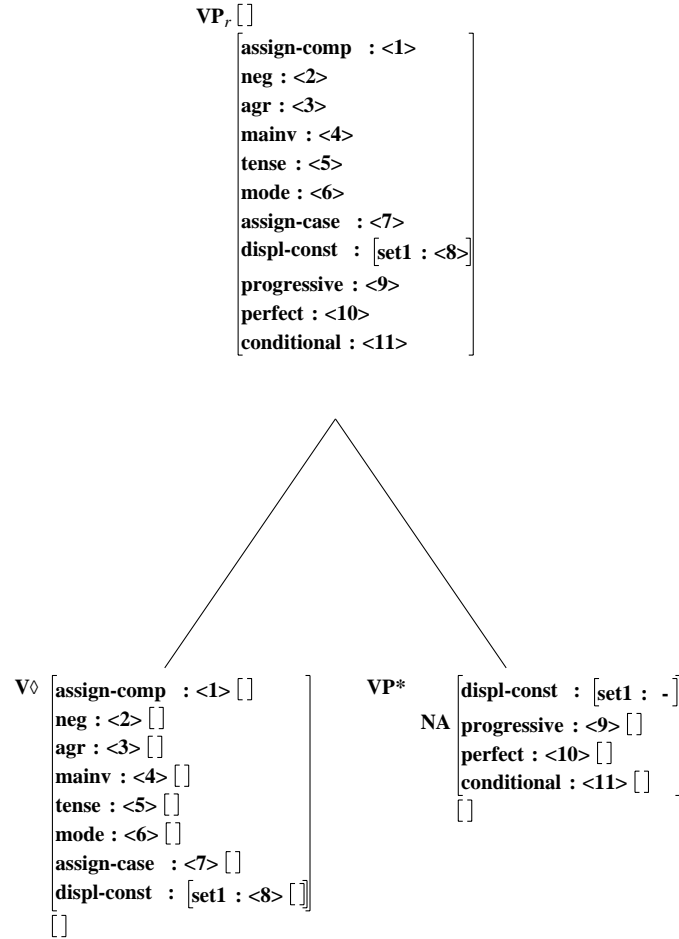
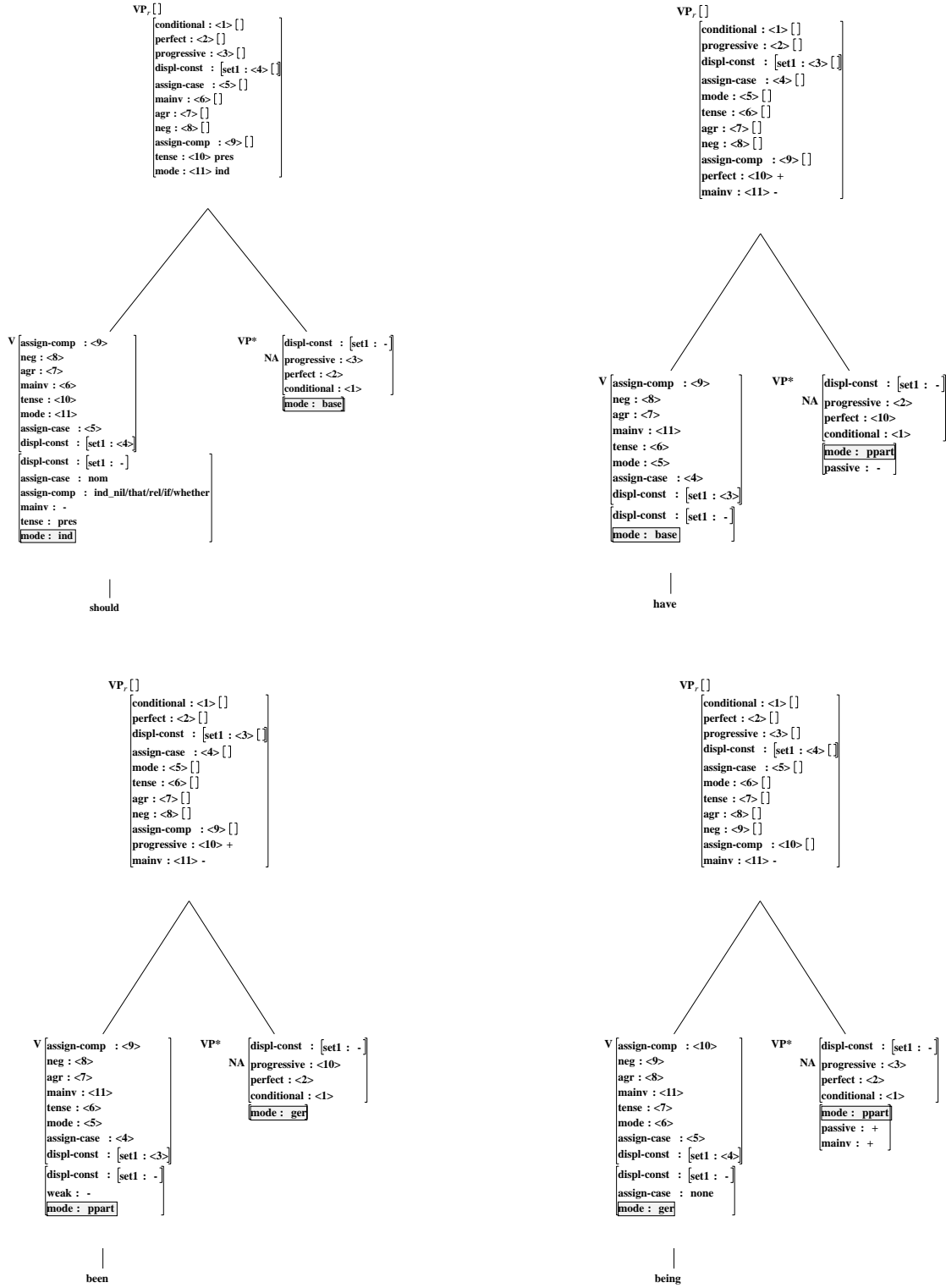


Figure 20.1: Auxiliary verb tree for non-inverted sentences: βV_{vx}

onto for the sentence to be valid, this enforces the restrictions made by the auxiliary node. In addition to these feature values, each auxiliary also gives values to the anchoring node ($V\Diamond$), to be passed up the tree to the root VP (VP_r) node; there they will become the new features for the top VP node of the sentential tree. Another auxiliary may now adjoin on top of it, and so forth. These feature values thereby ensure the proper auxiliary sequencing. Figure 20.2 shows the auxiliary trees anchored by the four auxiliary verbs in sentence (333). Figure 20.3 shows the final tree created for this sentence.

The general English restriction that matrix clauses must have tense (or be imperatives) is enforced by requiring the top S-node of a sentence to have **<mode>=ind/imp** (indicative or imperative). Since only the indicative and imperative sentences have tense, non-tensed clauses are restricted to occurring in embedded environments.

Noun-verb contractions are labeled NVC in their part-of-speech field in the morphological database and then undergo special processing to split them apart into the noun and the reduced verb before parsing. The noun then selects its trees in the normal fashion. The contraction, say


 Figure 20.2: Auxiliary trees for *The music should have been being played* .

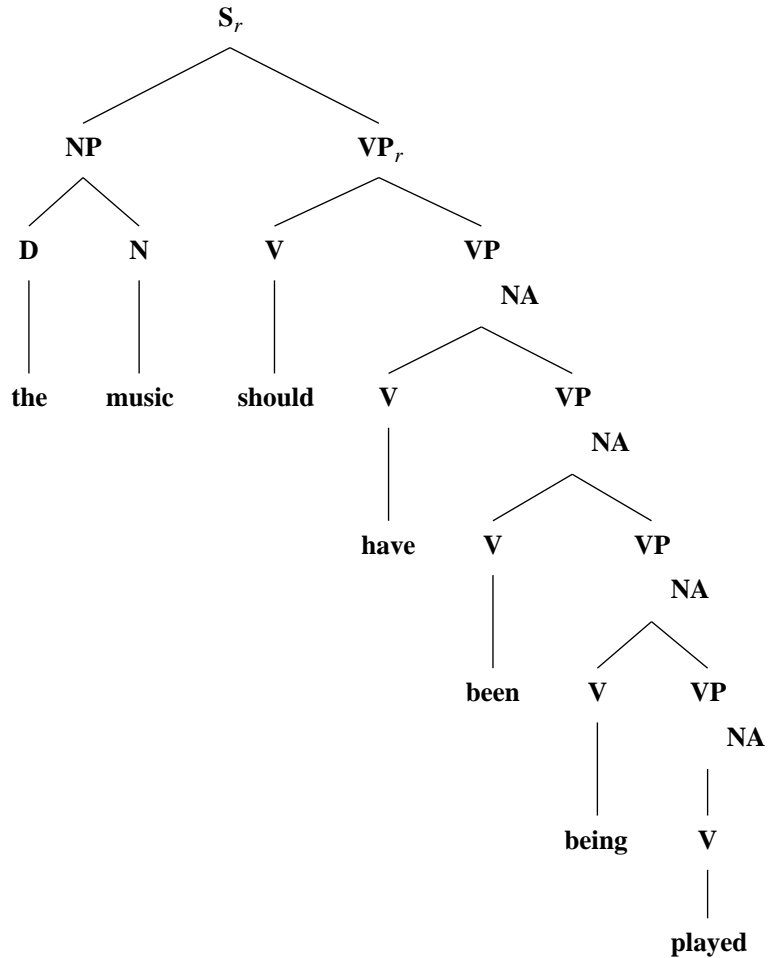


Figure 20.3: *The music should have been being played* .

'll or 'd, likewise selects the normal auxiliary verb tree, βVvx . However, since the contracted form, rather than the verb stem, is given in the morphology, the contracted form must also be listed as a separate syntactic entry. These entries have all the same features of the full form of the auxiliary verbs, with tense constraints coming from the morphological entry (e.g. *it's* is listed as IT 's NVC 3SG PRES). The ambiguous contractions 'd (*had/would*) and 's (*has/is*) behave like other ambiguous lexical items; there are simply multiple entries for those lexical items in the lexicon, each with different features. In the resulting parse, the contracted form is shown with features appropriate to the full auxiliary it represents.

20.2 Inverted Sentences

In inverted sentences, the two trees shown in Figure 20.4 adjoin to an S tree anchored by a main verb. The tree in Figure 20.4(a) is anchored by the auxiliary verb and adjoins to the S node, while the tree in Figure 20.4(b) is anchored by an empty element and adjoins at the VP

node. Figure 20.5 shows these trees (anchored by *will*) adjoined to the declarative transitive tree⁴ (anchored by main verb *buy*).

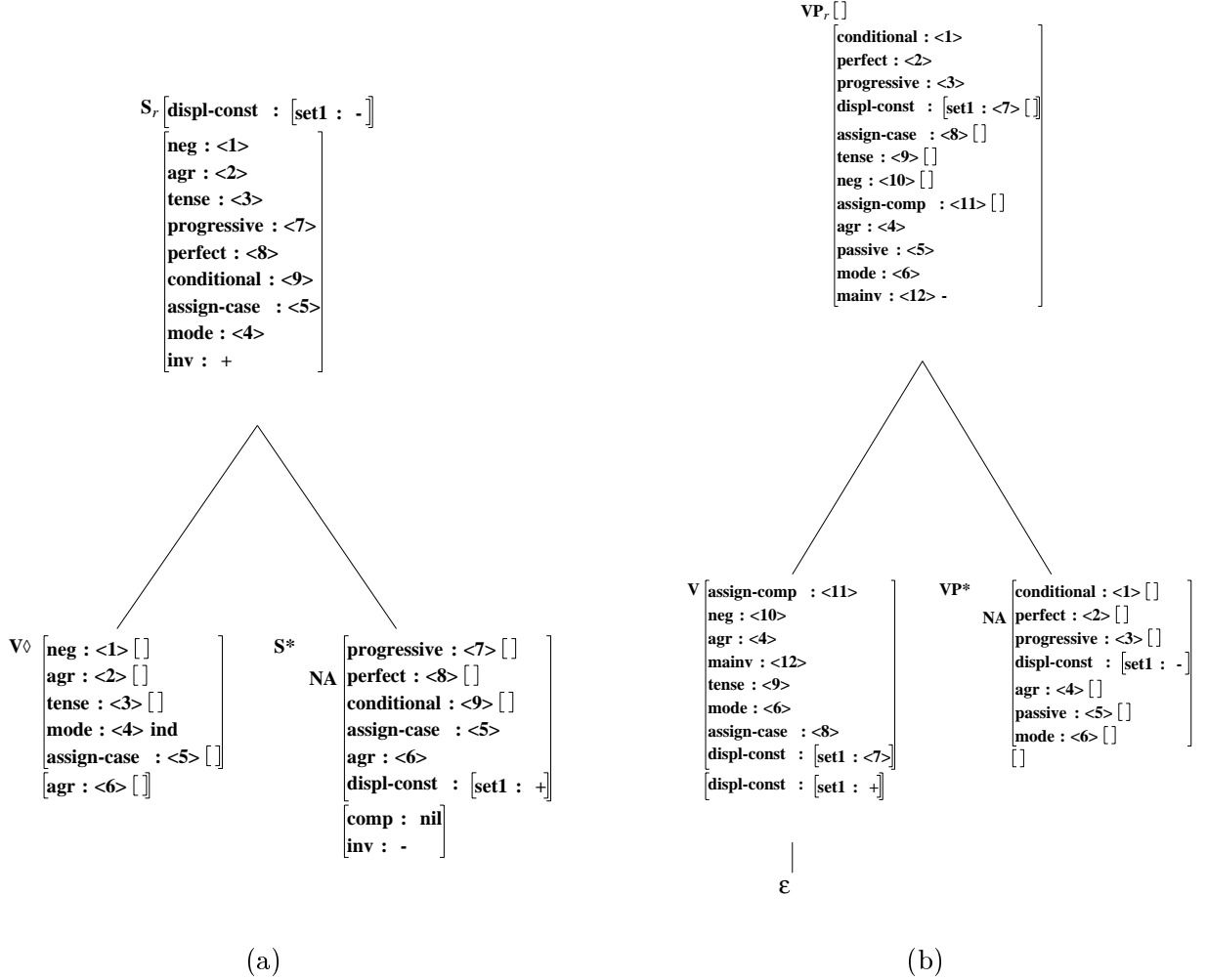


Figure 20.4: Trees for auxiliary verb inversion: βV_s (a) and βV_{vx} (b)

The feature **<displ-const>** ensures that both of the trees in Figure 20.4 must adjoin to an elementary tree whenever one of them does. For more discussion on this mechanism, which simulates tree local multi-component adjunction, see [Hockey and Srinivas, 1993]. The tree in Figure 20.4(b), anchored by ϵ , represents the originating position of the inverted auxiliary. Its adjunction blocks the **<assign-case>** values of the VP it dominates from being co-indexed with the **<case>** value of the subject. Since **<assign-case>** values from the VP are blocked, the **<case>** value of the subject can only be co-indexed with the **<assign-case>** value of the inverted auxiliary (Figure 20.4(a)). Consequently, the inverted auxiliary functions as the case-assigner for the subject in these inverted structures. This is in contrast to the situation in uninverted structures where the anchor of the highest (leftmost) VP assigns case to the subject

⁴The declarative transitive tree was seen in section 6.2.

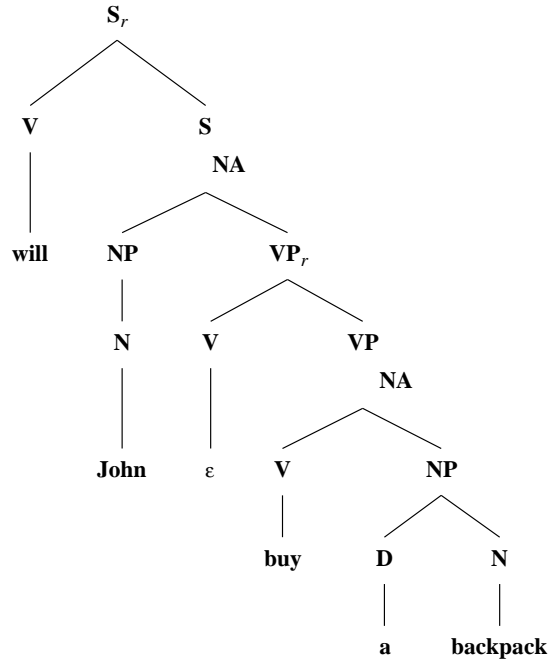


Figure 20.5: *will John buy a backpack ?*

(see section 4.4.3.2 for more on case assignment). The XTAG analysis is similar to GB accounts where the inverted auxiliary plus the ϵ -anchored tree are taken as representing I to C movement.

20.3 Do-Support

It is well-known that English requires a mechanism called ‘do-support’ for negated sentences and for inverted yes-no questions without auxiliaries.

(334) John does not want a car .

(335) *John not wants a car .

(336) John will not want a car .

(337) Do you want to leave home ?

(338) *want you to leave home ?

(339) will you want to leave home ?

20.3.1 In negated sentences

The GB analysis of *do*-support in negated sentences hinges on the separation of the INFL and VP nodes (see [Chomsky, 1965], [Jackendoff, 1972] and [Chomsky, 1986]). The claim is that the presence of the negative morpheme blocks the main verb from getting tense from the INFL node, thereby forcing the addition of a verbal lexeme to carry the inflectional elements. If an auxiliary verb is present, then it carries tense, but if not, periphrastic or ‘dummy’, *do* is required. This seems to indicate that *do* and other auxiliary verbs would not co-occur, and indeed this is the case (see sentences (340)-(341)). Auxiliary *do* is allowed in English when no negative morpheme is present, but this usage is marked as emphatic. Emphatic *do* is also not allowed to co-occur with auxiliary verbs (sentences (342)-(345)).

(340) *We will have *do* bought a sleeping bag .

(341) *We *do* will have bought a sleeping bag .

(342) You **do** have a backpack, don’t you ?

(343) I **do** want to go !

(344) *You **do** can have a backpack, don’t you ?

(345) *I **did** have had a backpack !

At present, the XTAG grammar does not have analyses for emphatic *do*.

In the XTAG grammar, *do* is prevented from co-occurring with other auxiliary verbs by a requirement that it adjoin only onto main verbs ($\langle \mathbf{mainv} \rangle = +$). It has indicative mode, so no other auxiliaries can adjoin above it.⁵ The lexical item *not* is only allowed to adjoin onto a non-indicative (and therefore non-tensed) verb. Since all matrix clauses must be indicative (or imperative), a negated sentence will fail unless an auxiliary verb, either *do* or another auxiliary, adjoins somewhere above the negative morpheme, *not*. In addition to forcing adjunction of an auxiliary, this analysis of *not* allows it freedom to move around in the auxiliaries, as seen in the sentences (346)-(349).

(346) John will have had a backpack .

(347) *John not will have had a backpack .

(348) John will not have had a backpack .

(349) John will have not had a backpack .

⁵Earlier, we said that indicative mode carries tense with it. Since only the topmost auxiliary carries the tense, any subsequent verbs must **not** have indicative mode.

20.3.2 In inverted yes/no questions

In inverted yes/no questions, *do* is required if there is no auxiliary verb to invert, as seen in sentences (337)-(339), replicated here as (350)-(352).

(350) do you want to leave home ?

(351) *want you to leave home ?

(352) will you want to leave home ?

(353) *do you will want to leave home ?

In English, unlike other Germanic languages, the main verb cannot move to the beginning of a clause, with the exception of main verb *be*.⁶ In a GB account of inverted yes/no questions, the tense feature is said to be in C^0 at the front of the sentence. Since main verbs cannot move, they cannot pick up the tense feature, and do-support is again required if there is no auxiliary verb to perform the role. Sentence (353) shows that *do* does not interact with other auxiliary verbs, even when in the inverted position.

In XTAG, trees anchored by a main verb that lacks tense are required to have an auxiliary verb adjoin onto them, whether at the VP node to form a declarative sentence, or at the S node to form an inverted question. *Do* selects the inverted auxiliary trees given in Figure 20.4, just as other auxiliaries do, so it is available to adjoin onto a tree at the S node to form a yes/no question. The mechanism described in section 20.3.1 prohibits *do* from co-occurring with other auxiliary verbs, even in the inverted position.

20.4 Infinitives

The infinitive *to* is considered an auxiliary verb in the XTAG system, and selects the auxiliary tree in Figure 20.1. *To*, like *do*, does not interact with the other auxiliary verbs, adjoining only to main verb base forms, and carrying infinitive mode. It is used in embedded clauses, both with and without a complementizer, as in sentences (354)-(356). Since it cannot be inverted, it simply does not select the trees in Figure 20.4.

(354) John wants to have a backpack .

(355) John wants Mary to have a backpack .

(356) John wants for Mary to have a backpack .

The usage of infinitival *to* interacts closely with the distribution of null subjects (PRO), and is described in more detail in section 8.5.

⁶The inversion of main verb *have* in British English was previously noted.

20.5 Semi-Auxiliaries

Under the category of semi-auxiliaries, we have placed several verbs that do not seem to closely follow the behavior of auxiliaries. One of these auxiliaries, *dare*, mainly behaves as a modal and selects for the base form of the verb. The other semi-auxiliaries all select for the infinitival form of the verb. Examples of this second type of semi-auxiliary are *used to*, *ought to*, *get to*, *have to*, and *BE to*.

20.5.1 Marginal Modal *dare*

The auxiliary *dare* is unique among modals in that it both allows DO-support and exhibits a past tense form. It clearly falls in modal position since no other auxiliary (except *do*) may precede it in linear order⁷. Examples appear below.

(357) she **dare** not have been seen .

(358) she does not **dare** succeed .

(359) Jerry **dared** not look left or right .

(360) only models **dare** wear such extraordinary outfits .

(361) **dare** Dale tell her the secret ?

(362) *Louise had dared not tell a soul .

As mentioned above, auxiliaries are prevented from having DO-support within the XTAG system. To allow for DO-support in this case, we had to create a lexical entry for *dare* that allowed it to have the feature **mainv+** and to have **base** mode (this measure is what also allows *dare* to occur in double-modal sequences). A second lexical entry was added to handle the regular modal occurrence of *dare*. Additionally, all other modals are classified as being present tense, while *dare* has both present and past forms. To handle this behavior, *dare* was given similar features to the other modals in the morphology minus the specification for tense.

20.5.2 Other semi-auxiliaries

The other semi-auxiliaries all select for the infinitival form of the verb. Many of these auxiliaries allow for DO-support and can appear in both base and past participle forms, in addition to being able to stand alone (indicative mode). Examples of this type appear below.

(363) Alex **used** to attend karate workshops .

(364) Angelina might have **used** to believe in fate .

(365) Rich did not **used** to want to be a physical therapist .

⁷Some speakers accept *dare* preceded by a modal, as in *I might dare finish this report today*. In the XTAG analysis, this particular double modal usage is accounted for. Other cases of double modal occurrence exist in some dialects of American English, although these are not accounted for in the system, as was mentioned earlier.

(366) Mick might not **have** to play the game tonight .

(367) Singer **had** to have been there .

(368) Heather has **got** to finish that project before she goes insane .

The auxiliaries *ought to* and *BE to* may not be preceded by any other auxiliary.

(369) Biff **ought** to have been working harder .

(370) *Carson does **ought** to have been working harder .

(371) the party **is** to take place this evening .

(372) *the party had **been** to take place this evening .

The trickiest element in this group of auxiliaries is *used to*. While the other verbs behave according to standard inflection for auxiliaries, *used to* has the same form whether it is in mode base, past participle, or indicative forms. The only connection *used to* maintains with the infinitival form *use* is that occasionally, the bare form *use* will appear with DO-support. Since the three modes mentioned above are mutually exclusive in terms of both the morphology and the lexicon, *used* has three entries in each.

20.5.3 Other Issues

There is a lingering problem with the auxiliaries that stems from the fact that there currently is no way to distinguish between the main verb and auxiliary verb behaviors for a given letter string within the morphology. This situation results in many unacceptable sentences being successfully parsed by the system. Examples of the unacceptable sentences are given below.

(373) the miller **cans** tell a good story . (vs the farmer **cans** peaches in July .)

(374) David **wills** have finished by noon . (vs the old man **wills** his fortune to me .)

(375) Sarah **needs** not leave . (vs Sarah **needs** to leave .)

(376) Jennifer **dares** not be seen . (vs the young woman **dares** him to do the stunt .)

(377) Lila **does use** to like beans . (vs Lila **does use** her new cookware .)

Chapter 21

Conjunction

21.1 Introduction

The XTAG system can handle sentences with conjunction of two constituents of the same syntactic category. The coordinating conjunctions which select the conjunction trees are *and*, *or* and *but*.¹ There are also multi-word conjunction trees, anchored by *either-or*, *neither-nor*, *both-and*, and *as well as*. There are eight syntactic categories that can be coordinated, and in each case an auxiliary tree is used to implement the conjunction. These eight categories can be considered as four different cases, as described in the following sections. In all cases the two constituents are required to be of the same syntactic category, but there may also be some additional constraints, as described below.

21.2 Adjective, Adverb, Preposition and PP Conjunction

Each of these four categories has an auxiliary tree that is used for conjunction of two constituents of that category. The auxiliary tree adjoins into the left-hand-side component, and the right-hand-side component substitutes into the auxiliary tree.

Figure 21.1(a) shows the auxiliary tree for adjective conjunction, and is used, for example, in the derivation of the parse tree for the noun phrase *the dark and dreary day*, as shown in Figure 21.1(b). The auxiliary tree adjoins onto the node for the left adjective, and the right adjective substitutes into the right hand side node of the auxiliary tree. The analysis for adverb, preposition and PP conjunction is exactly the same and there is a corresponding auxiliary tree for each of these that is identical to that of Figure 21.1(a) except, of course, for the node labels.

21.3 Noun Phrase and Noun Conjunction

The tree for NP conjunction, shown in Figure 21.2(a), has the same basic analysis as in the previous section except that the **<wh>** and **<case>** features are used to force the two noun phrases to have the same **<wh>** and **<case>** values. This allows, for example, *he and she wrote the book together* while disallowing **he and her wrote the book together*. Agreement is

¹We believe that the restriction of *but* to conjoining only two items is a pragmatic one, and our grammars accepts sequences of any number of elements conjoined by *but*.

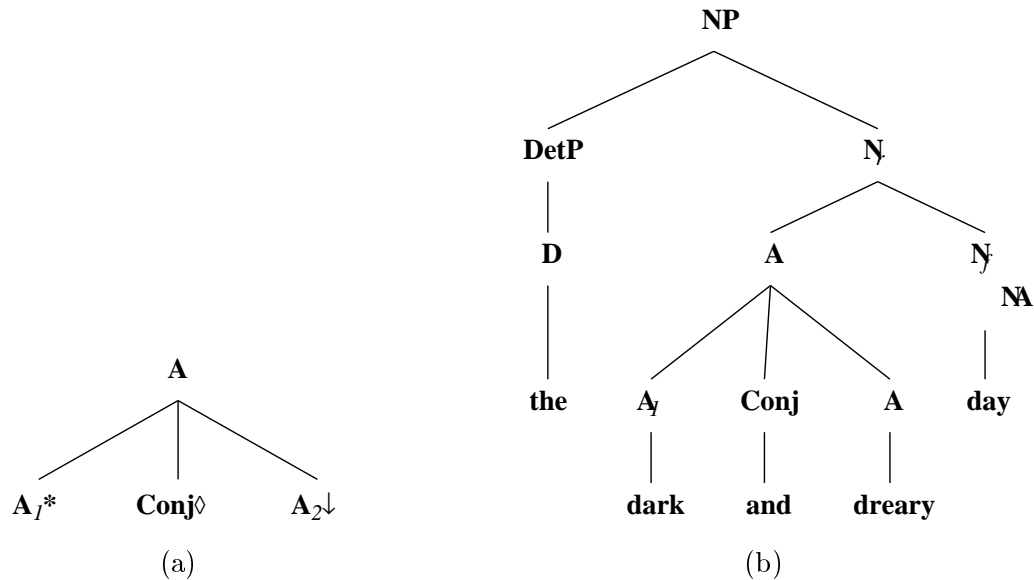


Figure 21.1: Tree for adjective conjunction: $\beta a1CONJa2$ and a resulting parse tree

lexicalized, since the various conjunctions behave differently. With *and*, the root $\langle \mathbf{agr\ num} \rangle$ value is $\langle \mathbf{plural} \rangle$, no matter what the number of the two conjuncts. With *or*, however, the root $\langle \mathbf{agr\ num} \rangle$ is co-indexed with the $\langle \mathbf{agr\ num} \rangle$ feature of the right conjunct. This ensures that the entire conjunct will bear the number of both conjuncts if they agree (Figure 21.2(b)), or of the most “recent” one if they differ (*Either the boys or John is going to help you.*). There is no rule per se on what the agreement should be here, but people tend to make the verb agree with the last conjunct (cf. Quirk et. al [1985], section 10.41 for discussion). The tree for N conjunction is identical to that for the NP tree except for the node labels. (The multi-word conjunctions do not select the N conjunction tree - **the both dogs and cats*).

21.4 Determiner Conjunction

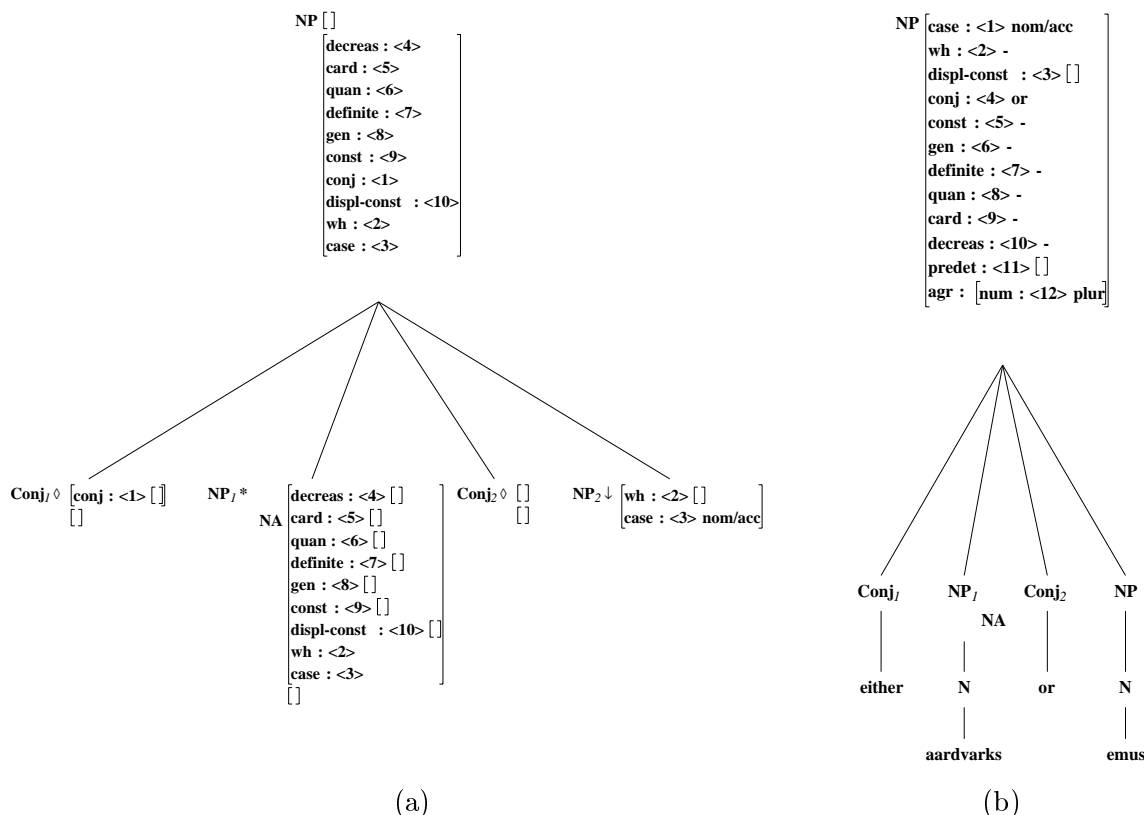
In determiner coordination, all of the determiner feature values are taken from the left determiner, and the only requirement is that the $\langle \mathbf{wh} \rangle$ feature is the same, while the other features, such as $\langle \mathbf{card} \rangle$, are unconstrained. For example, *which and what* and *all but one* are both acceptable determiner conjunctions, but **which and all* is not.

(378) how many and which people camp frequently ?

(379) *some or which people enjoy nature .

21.5 Sentential Conjunction

The tree for sentential conjunction, shown in Figure 21.4, is based on the same analysis as the conjunctions in the previous two sections, with a slight difference in features. The $\langle \mathbf{mode} \rangle$

Figure 21.2: Tree for NP conjunction: β CONJnx1CONJnx2 and a resulting parse tree

feature² is used to constrain the two sentences being conjoined to have the same mode so that *the day is dark and the phone never rang* is acceptable, but **the day dark and the phone never rang* is not. Similarly, the two sentences must agree in their **<wh>**, **<comp>** and **<extracted>** features. Co-indexation of the **<comp>** feature ensures that either both conjuncts have the same complementizer, or there is a single complementizer adjoined to the complete conjoined S. The **<assign-comp>** feature³ feature is used to allow conjunction of infinitival sentences, such as *to read and to sleep is a good life*.

21.6 Comma as a conjunction

We treat comma as a conjunction in conjoined lists. It anchors the same trees as the lexical conjunctions, but is considerably more restricted in how it combines with them. The trees anchored by commas are prohibited from adjoining to anything but another comma conjoined element or a non-coordinate element. (All scope possibilities are allowed for elements coordinated with lexical conjunctions.) Thus, structures such as Tree 21.5(a) are permitted, with each element stacking sequentially on top of the first element of the conjunct, while structures such as Tree 21.5(b) are blocked.

²See section 8.3 for an explanation of the **<mode>** feature.

³See section 8.5 for an explanation of the **<assign-comp>** feature.

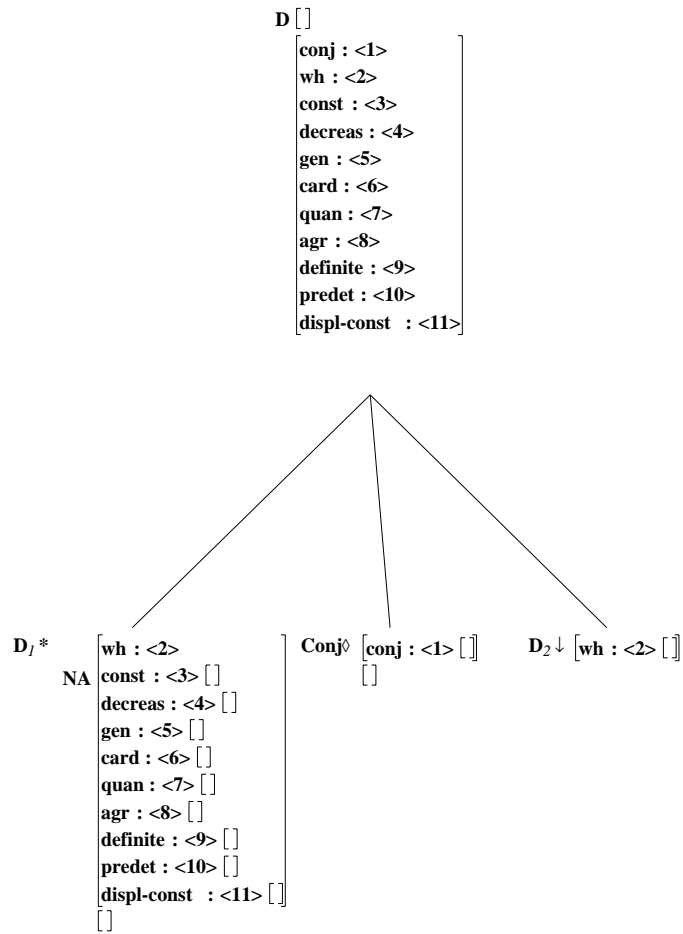


Figure 21.3: Tree for determiner conjunction: $\beta d1CONJd2.ps$

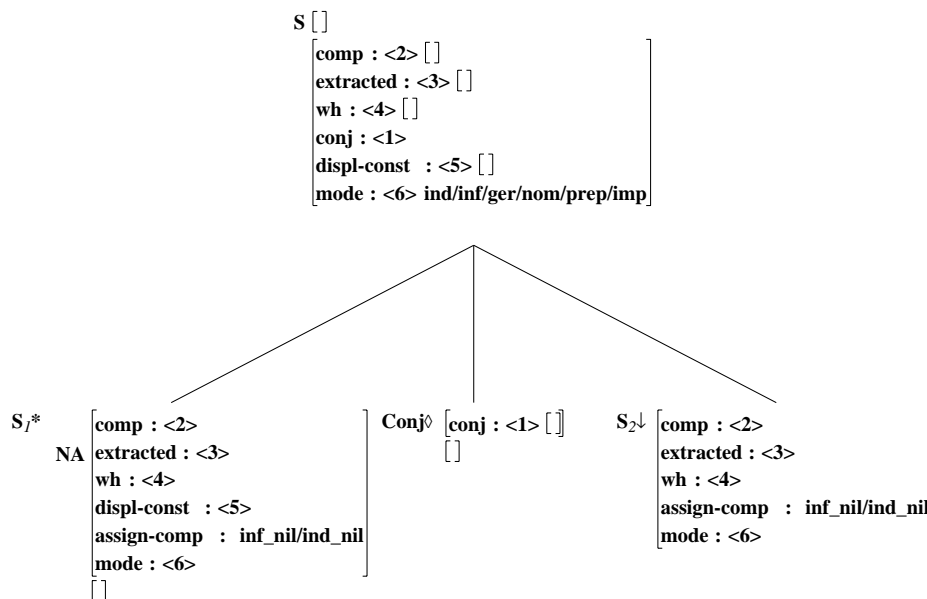
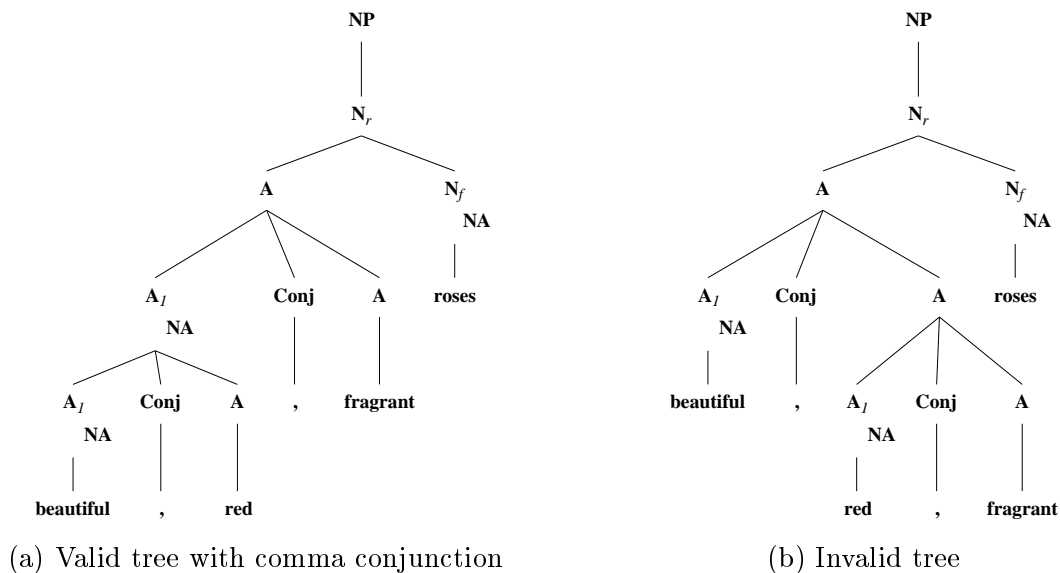

 Figure 21.4: Tree for sentential conjunction: $\beta s1CONJs2$


Figure 21.5:

This is accomplished by using the $\langle \text{conj} \rangle$ feature, which has the values **and/or/but** and **comma** to differentiate the lexical conjunctions from commas. The $\langle \text{conj} \rangle$ values for a comma-anchored tree and *and*-anchored tree are shown in Figure 21.6. The feature $\langle \text{conj} \rangle = \text{comma/none}$ on A_1 in (a) only allows comma conjoined or non-conjoined elements as the left-adjunct, and $\langle \text{conj} \rangle = \text{none}$ on A in (a) allows only a non-conjoined element as the right

conjunct. We also need the feature $\langle \text{conj} \rangle = \text{and/or/but/none}$ on the right conjunct of the trees anchored by lexical conjunctions like (b), to block comma-conjoined elements from substituting there. Without this restriction, we would get multiple parses of the NP in Tree 21.5; with the restrictions we only get the derivation with the correct scoping, shown as (a).

Since comma-conjoined lists can appear without a lexical conjunction between the final two elements, as shown in example (380), we cannot force all comma-conjoined sequences to end with a lexical conjunction.

- (380) So it is too with many other spirits which we all know: the spirit of Nazism or Communism, school spirit , the spirit of a street corner gang or a football team, the spirit of Rotary or the Ku Klux Klan. [Brown cd01]

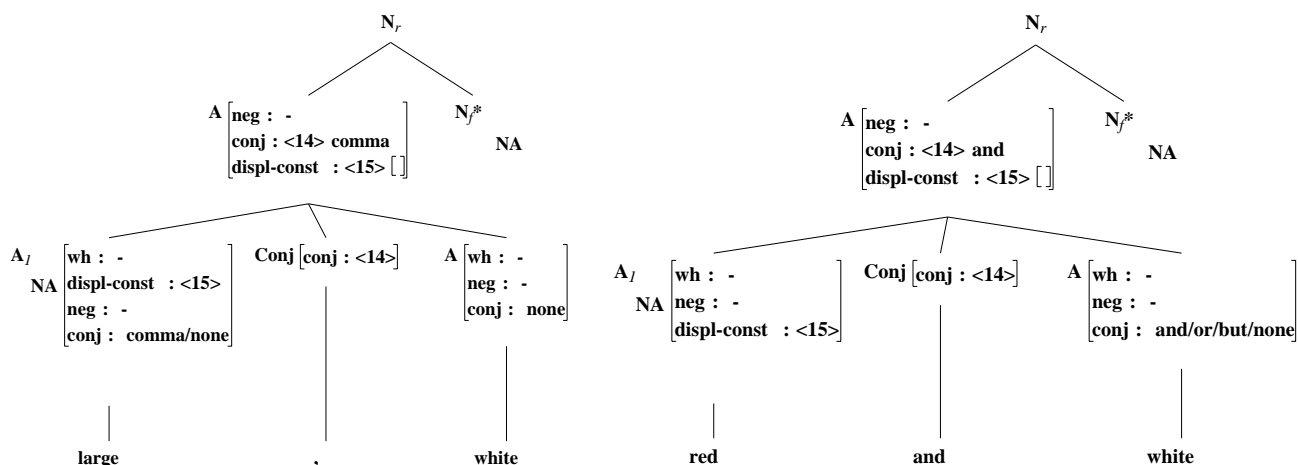


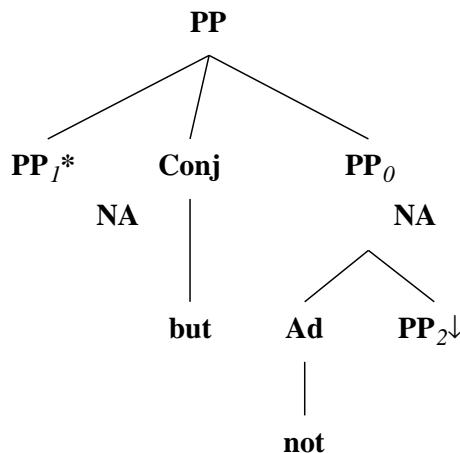
Figure 21.6: $\beta\text{a1CONJa2}$ (a) anchored by comma and (b) anchored by *and*

21.7 *But-not, not-but, and-not* and ϵ -not

We are analyzing conjoined structures such as *The women but not the men* with a multi-anchor conjunction tree anchored by the conjunction plus the adverb *not*. The alternative is to allow *not* to adjoin to any constituent. However, this is the only construction where *not* can freely occur onto a constituent other than a VP or adjective (cf. βNEGvx and βNEGa trees). It can also adjoin to some determiners, as discussed in Section 18. We want to allow sentences like (381) and rule out those like (382). The tree for the good example is shown in Figure 21.7. There are similar trees for *and-not* and ϵ -not, where ϵ is interpretable as either *and* or *but*, and a tree with *not* on the first conjunct for *not-but*.

- (381) Beth grows basil in the house (but) not in the garden .

- (382) *Beth grows basil (but) not in the garden .

Figure 21.7: Tree for conjunction with but-not: $\beta_{px1}CONJARB_{px2}$

Although these constructions sound a bit odd when the two conjuncts do not have the same number, they are sometimes possible. The agreement information for such NPs is always that of the non-negated conjunct: *his sons, and not Bill, are in charge of doing the laundry* or *not Bill, but his sons, are in charge of doing the laundry* (Some people insist on having the commas here, but they are frequently absent in corpus data.) The agreement feature from the non-negated conjunct is passed to the root NP, as shown in Figure 21.8. Aside from agreement, these constructions behave just like their non-negated counterparts.

21.8 *To* as a Conjunction

To can be used as a conjunction for adjectives (Fig. 21.9) and determiners, when they denote points on a scale:

(383) two to three degrees

(384) high to very high temperatures

As far as we can tell, when the conjuncts are determiners they must be cardinal.

21.9 Predicative Coordination

This section describes the method for predicative coordination (including VP coordination of various kinds) used in XTAG. The description is derived from work described in ([Sarkar and Joshi, 1996]). It is important to say that this implementation of predicative coordination is not part of the XTAG release at the moment due massive parsing ambiguities. This is partly because of the current implementation and also the inherent ambiguities due to VP coordination that cause a combinatorial explosion for the parser. We are trying to remedy both of these limitations using a probability model for coordination attachments which will be included as part of a later XTAG release.

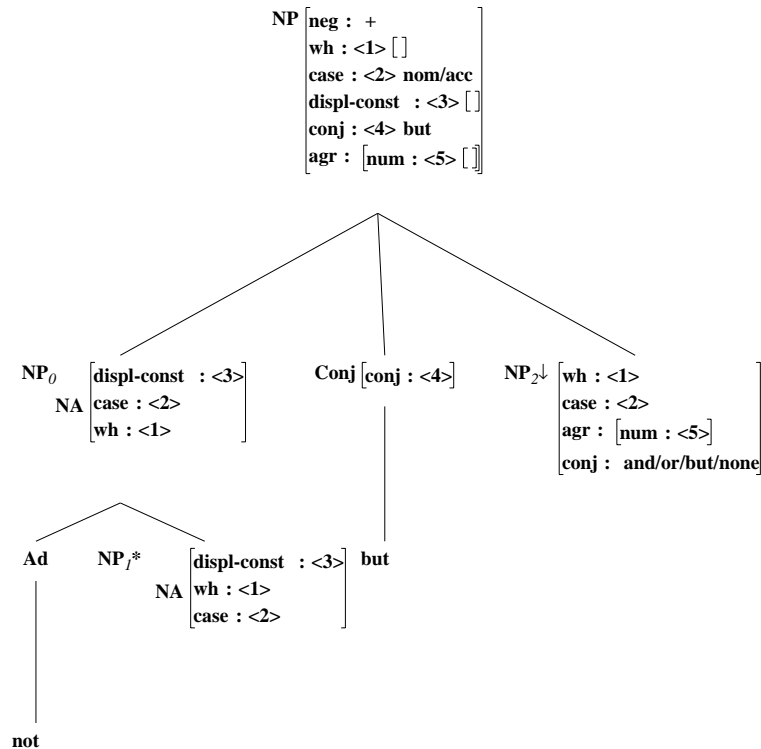


Figure 21.8: Tree for conjunction with not-but: $\beta\text{ARBnx1CONJnx2}$

This extended domain of locality in a lexicalized Tree Adjoining Grammar causes problems when we consider the coordination of such predicates. Consider (385) for instance, the NP *the beans that I bought from Alice* in the Right-Node Raising (RNR) construction has to be shared by the two elementary trees (which are anchored by *cooked* and *ate* respectively).

(385) (((Harry cooked) and (Mary ate)) the beans that I bought from Alice)

We use the standard notion of coordination which is shown in Figure 21.10 which maps two constituents of *like type*, but with different interpretations, into a constituent of the same type.

We add a new operation to the LTAG formalism (in addition to substitution and adjunction) called *conjoin* (later we discuss an alternative which replaces this operation by the traditional operations of substitution and adjunction). While substitution and adjunction take two trees to give a derived tree, *conjoin* takes three trees and composes them to give a derived tree. One of the trees is always the tree obtained by specializing the schema in Figure 21.10 for a particular category. The tree obtained will be a lexicalized tree, with the lexical anchor as the conjunction: *and*, *but*, etc.

The conjoin operation then creates a *contraction* between nodes in the contraction sets of the trees being coordinated. The term *contraction* is taken from the graph-theoretic notion of edge contraction. In a graph, when an edge joining two vertices is contracted, the nodes are merged and the new vertex retains edges to the union of the neighbors of the merged vertices. The conjoin operation supplies a new edge between each corresponding node in the contraction set and then contracts that edge.

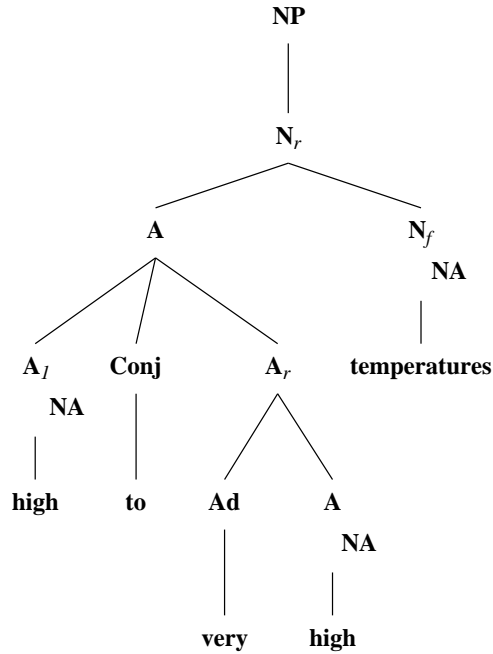
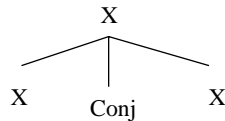

 Figure 21.9: Example of conjunction with *to*


Figure 21.10: Coordination schema

For example, applying *conjoin* to the trees *Conj(and)*, $\alpha(eats)$ and $\alpha(drinks)$ gives us the derivation tree and derived structure for the constituent in 386 shown in Figure 21.11.

(386) ... eats cookies and drinks beer

Another way of viewing the *conjoin* operation is as the construction of an auxiliary structure from an elementary tree. For example, from the elementary tree $\alpha(drinks)$, the *conjoin* operation would create the auxiliary structure $\beta(drinks)$ shown in Figure 21.12. The adjunction operation would now be responsible for creating contractions between nodes in the contraction sets of the two trees supplied to it. Such an approach is attractive for two reasons. First, it uses only the traditional operations of substitution and adjunction. Secondly, it treats *conj X* as a kind of “modifier” on the left conjunct *X*. This approach reduces some of the parsing ambiguities introduced by the predicative coordination trees and forms the basis of the XTAG implementation.

More information about predicative coordination can be found in ([Sarkar and Joshi, 1996]), including an extension to handle gapping constructions.

21.10 Pseudo-coordination

The XTAG grammar does handle one sort of verb pseudo-coordination. Semi-idiomatic phrases such as 'try and' and 'up and' (as in 'they might try and come today') are handled as multi-anchor modifiers rather than as true coordination. These items adjoin to a V node, using the β VCONJ_v tree. This tree adjoins only to verbs in their base morphological (non-inflected) form. The verb anchor of the β VCONJ_v must also be in its base form, as shown in examples (387)-(389). This blocks 3rd-person singular derivations, which are the only person morphologically marked in the present, except when an auxiliary verb is present or the verb is in the infinitive.

(387) *He tried and came yesterday.

(388) They try and exercise three times a week.

(389) He wants to try and sell the puppies.

Chapter 22

Comparatives

22.1 Introduction

Comparatives in English can manifest themselves in many ways, acting on many different grammatical categories and often involving ellipsis. A distinction must be made at the outset between two very different sorts of comparatives—those which make a comparison between two propositions and those which compare the extent to which an entity has one property to a greater or lesser extent than another property. The former, which we will refer to as *propositional* comparatives, is exemplified in (390), while the latter, which we will call *metalinguistic* comparatives (following Hellan 1981), is seen in (391):

(390) Ronaldo is more angry than Romario.

(391) Ronaldo is more angry than upset.

In (390), the extent to which Ronaldo is angry is greater than the extent to which Romario is angry. Sentence (391) indicates that the extent to which Ronaldo is angry is greater than the extent to which he is upset.

Apart from certain of the elliptical cases, both kinds of comparatives can be handled straightforwardly in the XTAG system. Elliptical cases which are not presently covered include those exemplified by the following sentences, which would presumably be handled in the same way as other sorts of VP ellipsis would.

(392) Ronaldo is more angry than Romario is.

(393) Bill eats more broccoli than George eats.

(394) Bill eats more broccoli than George does.

We turn to the analysis of metalinguistic comparatives first.

22.2 Metalinguistic Comparatives

A metalinguistic comparison can be performed on basically all of the predication categories—adjectives, verb phrases, prepositional phrases, and nouns—as in the following examples:

- (395) The table is more long than wide. (AP)
 (396) Clark more makes the rules than follows them. (VP)
 (397) Calvin is more in the living room than in the kitchen. (PP)
 (398) That unidentified amphibian in the bush is more frog than toad, I would say. (NP)

At present, we only deal with the adjectival metalinguistic comparatives as in (395). The analysis given here for these can be easily extended to prepositional phrases and nominal comparatives of the metalinguistic sort, but, as with coordination in XTAG, verb phrases will prove more difficult.

Adjectival comparatives appear to distribute with simple adjectives, as in the following examples:

- (399) Herbert is more livid than angry.
 (400) Herbert is more livid and furious than angry.
 (401) The more innovative than conventional medication cured everyone in the sick ward.
 (402) The elephant, more wobbly than steady, fell from the circus ball.

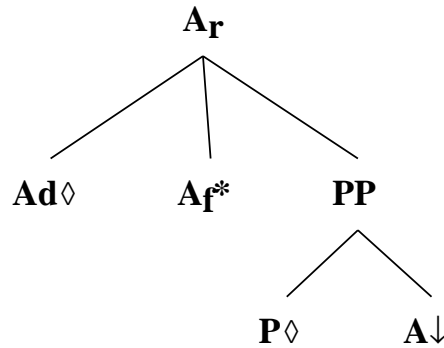


Figure 22.1: Tree for Metalinguistic Adjective Comparative: β ARBaPa

This patterning indicates that we can give these comparatives a tree that adjoins quite freely onto adjectives, as in Figure 22.1. This tree is anchored by *more/less - than*. To avoid grammatically incorrect comparisons such as *more brighter than dark*, the feature **compar** is used to block this tree from adjoining onto morphologically comparative adjectives. The foot node is **compar-**, while *brighter* and its comparative siblings are **compar+**¹. We also wish to block strings like *more brightest than dark*, which is accomplished with the feature **super**, indicating superlatives. This feature is negative at the foot node so that β ARBaPa cannot adjoin to superlatives like *niciest*, which are specified as **super+** from the morphology. Furthermore, the root node is **super+** so that β ARBaPa cannot adjoin onto itself and produce monstrosities such as (403):

¹The analysis given later for adjectival propositional comparatives produces aggregated **compar+** adjectives such as *more bright*, which will also be incompatible (as desired) with β ARBaPa.

(403) *Herbert is more less livid than angry than furious.

Thus, the use of the **super** feature is less to indicate superlativeness specifically, but rather to indicate that the subtree below a **super**+ node contains a full-fledged comparison. In the case of lexical superlatives, the comparison is against everything, implicitly.

A benefit of the multiple-anchor approach here is that we will never allow sentences such as (404), which would be permissible if we split the comparative component and the *than* component of metalinguistic comparatives into two separate trees.

(404) *Ronaldo is angrier than upset.

We also see another variety of adjectival comparatives of the form *more/less than X*, which indicates some property which is more or less extreme than the property *X*. In a sentence such as (405), some property is being said to hold of Francis such that it is of a kind with *stupid* and that it exceeds *stupid* on some scale (intelligence, for example). Quirk et al. also note that these constructions remark on the inadequacy of the lexical item. Thus, in (404), it could be that *stupid* is a starting point from which the speaker makes an approximation for some property which the speaker feels is beyond the range of the English lexicon, but which expresses the supreme lack of intellect of the individual it is predicated of.

(405) Francis is more than stupid.

(406) Romario is more than just upset.

Taking our inspiration from β ARBaPa, we can handle these comparatives, which have the same distribution but contain an empty adjective, by using the tree shown in Figure 22.2.

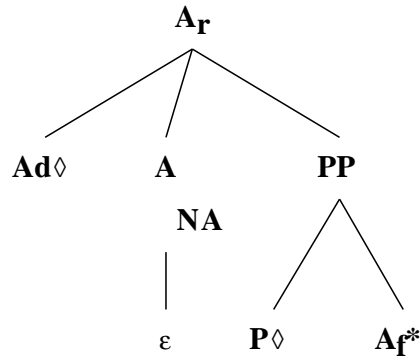


Figure 22.2: Tree for Adjective-Extreme Comparative: β ARBaPa

This sort of metalinguistic comparative also occurs with the verb phrase, prepositional phrase, and noun varieties.

(407) Clark more than makes the rules. (VP)

(408) Calvin's hands are more than near the cookie jar. (PP)

(409) That stuff on her face is more than mud. (NP)

Presumably, the analysis for these would parallel that for adjectives, though it has not yet been implemented.

22.3 Propositional Comparatives

22.3.1 Nominal Comparatives

Nominal comparatives are considered here to be those which compare the cardinality of two sets of entities denoted by nominal phrases. The following data lay out a basic distribution of these comparatives.

(410) More vikings than mongols eat spam.

(411) *More the vikings than mongols eat spam.

(412) Vikings eat less spaghetti than spam.

(413) More men that walk to the store than women who despise spam enjoyed the football game.

(414) More men than James like scotch on the rocks.

(415) Elmer knows fewer martians than rabbits.

Looking at these examples, we are tempted to produce a tree for this construction that is similar to βARBaPa . However, it is quite common for the *than* portion of these comparatives to be left out, as in the following sentences:

(416) More vikings eat spam.

(417) Mongols eat less spam.

Furthermore, *than NP* cannot occur without *more*. These facts indicate that we can and should build up nominal comparatives with two separate trees. The first, which allows a comparative adverb to adjoin to a noun, is given in Figure 22.3(a). The second is the noun-phrase modifying prepositional tree. The tree βCARBn is anchored by *more/less/fewer* and βCnxPnx is anchored by *than*. The feature **compar** is used to ensure that only one βCARBn tree can adjoin to any given noun—its foot node is **compar-** and the root node is **compar+**. All nouns are **compar-**, and the **compar** value is passed up through all trees which adjoin to N or NP. In order to ensure that we do not allow sentences like **Vikings than mongols eat spam*, the **compar** feature is used. The NP foot node of βCnxPnx is **compar+**; thus, βCnxPnx will adjoin only to NP's which have been already modified by βCARBn (and thereby comparativized). In this way, we capture sentences like (416) en route to deriving sentences like (410), in a principled and simple manner.

Further evidence for this approach comes from comparative clauses which are missing the noun phrase which is being compared against something, as in the following:

(418) The vikings ate more.²

²We ignore here the interpretation in which the comparison covers the eating event, focussing only on the one which the comparison involves the stuff being eaten.

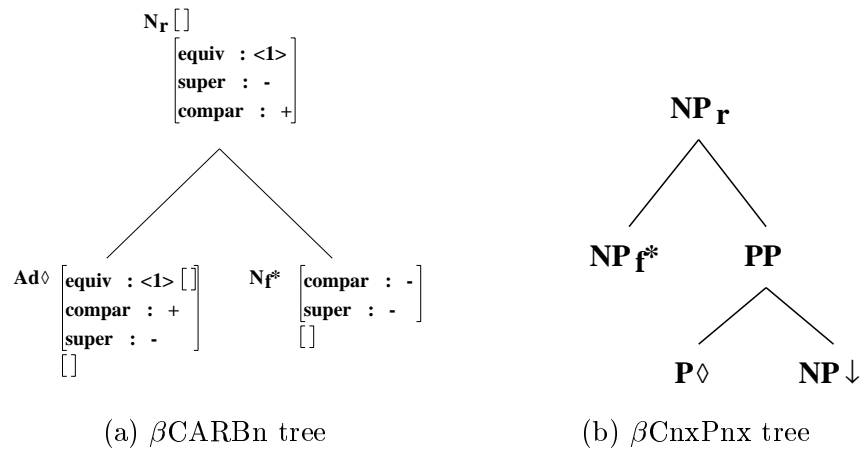


Figure 22.3: Nominal comparative trees

(419) The vikings ate more than a boar.³

Sometimes the missing noun refers to an entity or set available in the prior discourse, while at other times it is a reference to some anonymous, unspecified set. The former is exemplified in a mini-discourse such as the following:

Calvin: “The mongols ate spam.”

Hobbes: “The vikings ate more.”

The latter can be seen in the following example:

Calvin: “The vikings ate a a boar.”

Hobbes: “Indeed. But in fact, the vikings ate more than a boar.”

Since the lone comparatives *more/less/fewer* have the same basic distribution as noun phrases, the tree in Figure 22.4 is employed to capture this fact. The root node of α CARB is **compar**+. Not only does this accord with our intuitions about what the **compar** feature is supposed to indicate, it also permits β_{nxPnx} to adjoin, giving us strings such as *more than NP* for free.

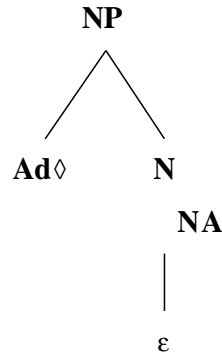


Figure 22.4: Tree for Lone Comparatives: α CARB

Thus, by splitting nominal comparatives into multiple trees, we make correct predictions about their distribution with a minimal number of simple trees. Furthermore, we now also get certain comparative coordinations for free, once we place the requirement that nouns and noun phrases must match for **compar** if they are to be coordinated. This yields strings such as the following:

(420) Julius eats more grapes and fewer boars than avocados.

(421) Were there more or less than fifty people (at the party)?

³This sentence differs from the metalinguistic comparison *That stuff on her face is more than mud* in that it involves a comment on the quantity and/or type of the compared NP, whereas the other expresses that the property denoted by the compared noun is an inadequate characterization of the thing being described.

The structures are given in Figure 22.5. Also, it will block strings like *more men and women than children* under the (impossible) interpretation that there are more men than children but the comparison of the quantity of women to children is not performed. Unfortunately, it will permit comparative clauses such as *more grapes and fewer than avocados* under the interpretation in which there are more grapes than avocados and fewer of some unspecified thing than avocados (see Figure 22.6).

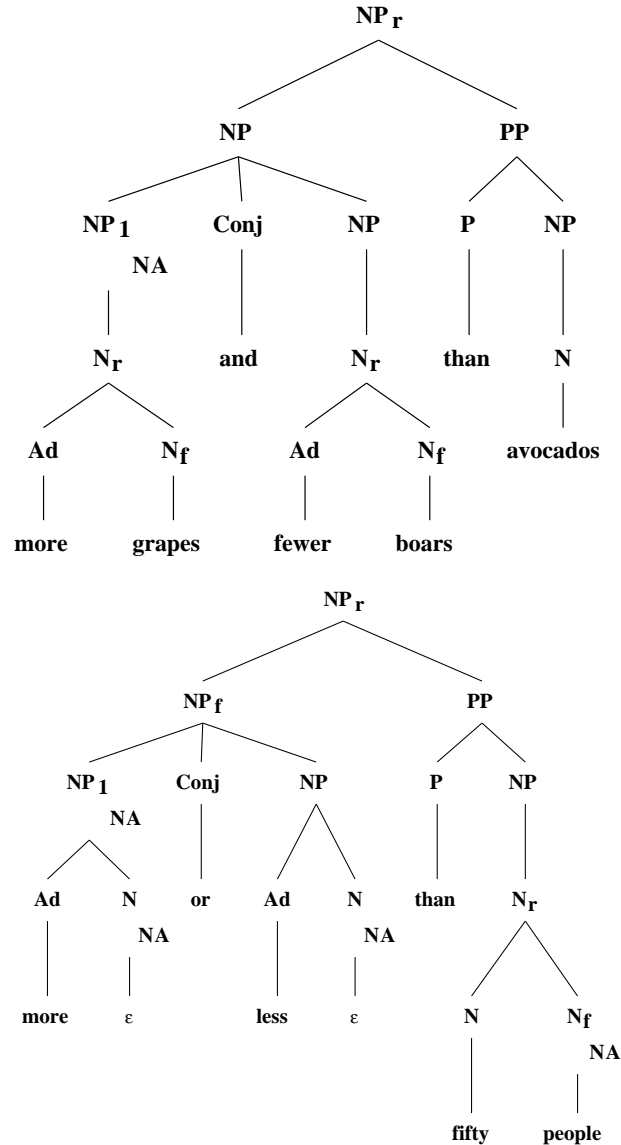


Figure 22.5: Comparative conjunctions.

One aspect of this analysis is that it handles the elliptical comparatives such as the following:

(422) Arnold kills more bad guys than Steven.

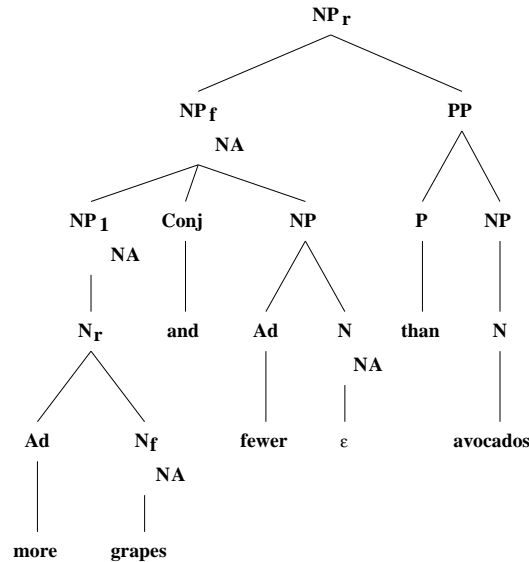


Figure 22.6: Comparative conjunctions.

In a sense, this is actually only simulating the ellipsis of these constructions indirectly. However, consider the following sentences:

(423) Arnold kills more bad guys than I do.

(424) Arnold kills more bad guys than I.

(425) Arnold kills more bad guys than me.

The first of these has a *pro*-verb phrase which has a nominative subject. If we totally drop the second verb phrase, we find that the second NP can be in either the nominative or the accusative case. Prescriptive grammars disallow accusative case, but it actually is more common to find accusative case—use of the nominative in conversation tends to sound rather stiff and unnatural. This accords with the present analysis in which the second noun phrase in these comparatives is the complement of *than* in $\beta_{\text{nx}}\text{P}_{\text{nx}}$, and receives its case-marking from *than*. This does mean that the grammar will not currently accept (424), and indeed such sentences will only be covered by an analysis which really deals with the ellipsis. Yet the fact that most speakers produce (425) indicates that some sort of restructuring has occurred that results in the kind of structure the present analysis offers.

There is yet another distributional fact which falls out of this analysis. When comparative or comparativized adjectives modify a noun phrase, they can stand alone or occur with a *than* phrase; furthermore, they are obligatory when a *than*-phrase is present.

(426) Hobbes is a better teacher.

(427) Hobbes is a better teacher than Bill.

(428) A more exquisite horse launched onto the racetrack.

(429) A more exquisite horse than Black Beauty launched onto the racetrack.

(430) *Hobbes is a teacher than Bill.

Comparative adjectives such as *better* come from the lexicon as **compar+**. By having trees such as βAn transmit the **compar** value of the A node to the root N node, we can signal to βCnxPnx that it may adjoin when a comparative adjective has adjoined. An example of such an adjunction is given in Figure 22.7. Of course, if no comparative element is present in the lower part of the noun phrase, βnxPnx will not be able to adjoin since nouns themselves are **compar-**. In order to capture the fact that a comparative element blocks further modification to N, βAn must only adjoin to N nodes which are **compar-** in their lower feature matrix.

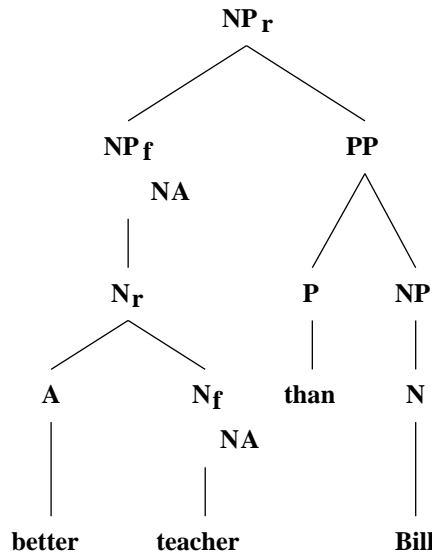


Figure 22.7: Adjunction of βnxPnx to NP modified by comparative adjective.

In order to obtain this result for phrases like *more exquisite horse*, we need to provide a way for *more* and *less* to modify adjectives without a *than*-clause as we have with βARBaPa . Actually, we need this ability independently for comparative adjectival phrases, as discussed in the next section.

22.3.2 Adjectival Comparatives

With nominal comparatives, we saw that a single analysis was amenable to both “pure” comparatives and elliptical comparatives. This is not possible for adjectival comparatives, as the following examples demonstrate:

(431) The dog is less patient.

(432) The dog is less patient than the cat.

(433) The dog is as patient.

(434) The dog is as patient as the cat.

(435) The less patient dog waited eagerly for its master.

(436) *The less patient than the cat dog waited eagerly for its master.

The last example shows that comparative adjectival phrases cannot distribute quite as freely as comparative nominals.

The analysis of elliptical comparative adjectives follows closely to that of comparative nominals. We build them up by first adjoining the comparative element to the A node, which then signals to the AP node, via the **compar** feature, that it may allow a *than*-clause to adjoin. The relevant trees are given in Figure 22.8. β CARBa is anchored by *more*, *less* and *as*, and β axPnx is anchored by both *than* and *as*.

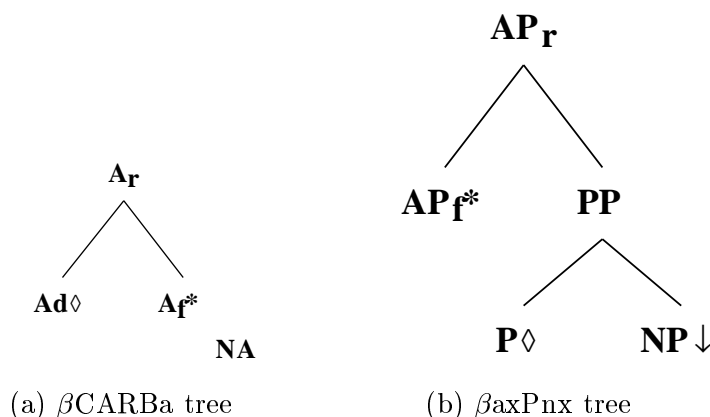


Figure 22.8: Elliptical adjectival comparative trees

The advantages of this analysis are many. We capture the distribution exhibited in the examples given in (431) - (436). With β CARBa, comparative elements may modify adjectives wherever they occur. However, *than* clauses for adjectives have a more restricted distribution which coincides nicely with the distribution of AP's in the XTAG grammar. Thus, by making them adjoin to AP rather than A, ill-formed sentences like (436) are not allowed.

There are two further advantages to this analysis. One is that β CARBa interacts with β nxPnx to produce sequences like *more exquisite horse than Black Beauty*, a result alluded to at the end of Section 22.3.1. We achieve this by ensuring that the comparativeness of an adjective is controlled by a comparative adverb which adjoins to it. A sample derivation is given in Figure 22.9. The second advantage is that we get sentences such as (437) for free.

(437) Hobbes is better than Bill.

Since *better* comes from the lexicon as **compar+** and this value is passed up to the AP node, β axPnx can adjoin as desired, giving us the derivation given in Figure 22.10.

Notice that the root AP node of Figure 22.10 is **compar-**, so we are basically saying that strings such as *better than Bill* are not “comparative.” This accords with our use of the **compar** feature—a positive value for **compar** signals that the clause beneath it is to **be** compared

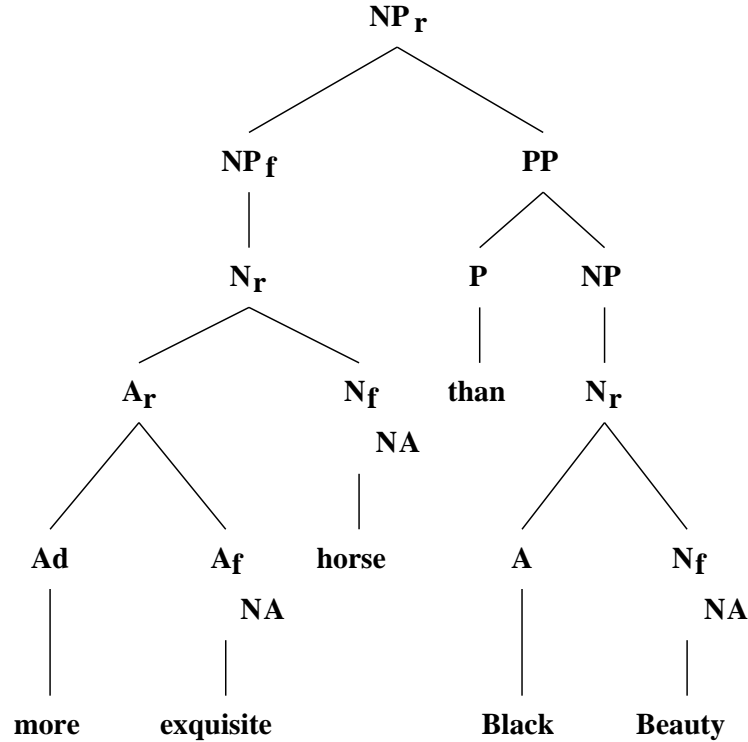


Figure 22.9: Comparativized adjective triggering βCuxPnx .

against something else. In the case of *better than Bill*, the comparison has been fulfilled, so we do not want it to signal for further comparisons. A nice result which follows is that βaxPnx cannot adjoin more than once to any given AP spine, and we have no need for the NA constraint on the tree's root node. Also, this treatment of the comparativeness of various strings proves important in getting the coordination of comparative constructions to work properly.

A note needs to be made about the analysis regarding the interaction of the equivalence comparative construction *as ... as* and the inequivalence comparative construction *more/less ... than*. In the grammar, *more*, *less*, and *as* all anchor βCARBa , and both *than* and *as* anchor βaxPnx . Without further modifications, this of course will give us sentences such as the following:

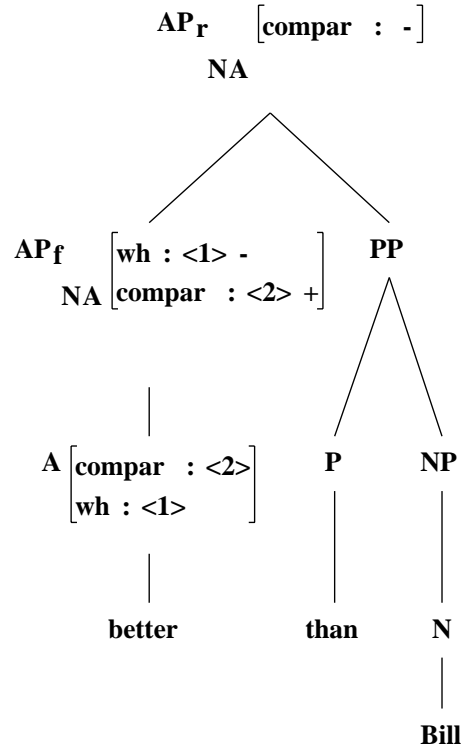
(438) *?Hobbes is as patient than Bill.

(439) *?Hobbes is more patient as Bill.

Such cases are blocked with the feature **equiv**: *more*, *less*, *fewer* and *than* are **equiv-** while *as* (in both adverbial and prepositional uses) is **equiv+**. The prepositional trees then require that their P node and the node to which they are adjoining match for **equiv**.

An interesting phenomena in which comparisons seem to be paired with an inappropriate *as/than*-clause is exhibited in (440) and (441).

(440) Hobbes is as patient or more patient than Bill.

Figure 22.10: Adjunction of $\beta_{axP_{nx}}$ to comparative adjective.

(441) Hobbes is more patient or as patient as Bill.

Though prescriptive grammars disfavor these sentences, these are perfectly acceptable. We can capture the fact that the *as/than*-clause shares the **equiv** value with the latter of the comparison phrases by passing the **equiv** value for the second element to the root of the coordination tree.

22.3.3 Adverbial Comparatives

The analysis of adverbial comparatives encouragingly parallels the analysis for nominal and elliptical adjectival comparatives—with, however, some interesting differences. Some examples of adverbial comparatives and their distribution are given in the following:

(442) Albert works more quickly.

(443) Albert works more quickly than Richard.

(444) Albert works more.

(445) *Albert more works.

(446) Albert works more than Richard.

(447) Hobbes eats his supper more quickly than Calvin.

(448) Hobbes more quickly eats his supper than Calvin.

(449) *Hobbes more quickly than Calvin eats his supper.

When *more* is used alone as an adverb, it must also occur after the verb phrase. Also, it appears that adverbs modified by *more* and *less* have the same distribution as when they are not modified. However, the *than* portion of an adverbial comparative is restricted to post verb phrase positions.

The first observation can be captured by having *more* and *less* select only $\beta vxARB$ from the set of adverb trees. Comparativization of adverbs looks very similar to that of other categories, and we follow this trend by giving the tree in Figure 22.11(a), which parallels the adjectival and nominal trees, for these instances. This handles the quite free distribution of adverbs which have been comparativized, while the tree in Figure 22.11(b), $\beta vxPnx$, allows the *than* portion of an adverbial comparative to occur only after the verb phrase, blocking examples such as (449).

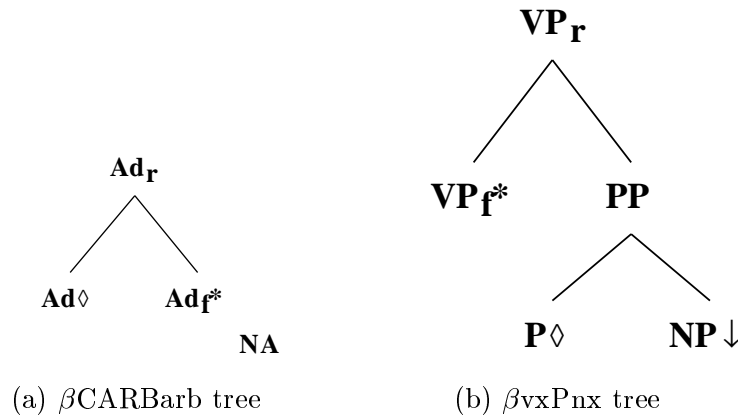


Figure 22.11: Adverbial comparative trees

The usage of the **compar** feature parallels that of the adjectives and nominals; however, trees which adjoin to VP are **compar-** on their root VP node. In this way, $\beta vxPnx$ anchored by *than* or *as* (which must adjoin to a **compar+** VP) can only adjoin immediately above a comparative or comparativized adverb. This avoids extra parses in which the comparative adverb adjoins at a VP node lower than the *than*-clause.

A final note is that *as* may anchor $\beta vxPnx$ non-comparatively, as in sentence (450). This means that there will be two parses for sentences such as (451).

(450) John works as a carpenter.

(451) John works as quickly as a carpenter.

This appears to be a legitimate ambiguity. One is that John works as quickly as a carpenter (works quickly), and the other is that John works quickly when he is acting as a carpenter (but maybe he is slow when he acting as a plumber).

22.4 Future Work

- Interaction with determiner sequencing (e.g., *several more men than women* but not **every more men than women*).
- Handle sentential complement comparisons (e.g., *Bill eats more pasta than Angus drinks beer*).
- Add partitives.
- Deal with constructions like *as many* and *as much*.
- Look at *so...as* construction.

Chapter 23

Punctuation Marks

Many parsers require that punctuation be stripped out of the input. Since punctuation is often optional, this sometimes has no effect. However, there are a number of constructions which must obligatorily contain punctuation and adding analyses of these to the grammar without the punctuation would lead to severe overgeneration. An especially common example is noun appositives. Without access to punctuation, one would have to allow every combinatorial possibility of NPs in noun sequences, which is clearly undesirable (especially since there is already unavoidable noun-noun compounding ambiguity). Aside from coverage issues, it is also preferable to take input “as is” and do as little editing as possible. With the addition of punctuation to the XTAG grammar, we need only do/assume the conversion of certain sequences of punctuation into the “British” order (this is discussed in more detail below in Section 23.2).

The XTAG POS tagger currently tags every punctuation mark as itself. These tags are all converted to the POS tag *Punct* before parsing. This allows us to treat the punctuation marks as a single POS class. They then have features which distinguish amongst them. Wherever possible we have the punctuation marks as anchors, to facilitate early filtering.

The full set of punctuation marks is separated into three classes: balanced, separating and terminal. The balanced punctuation marks are quotes and parentheses, separating are commas, dashes, semi-colons and colons, and terminal are periods, exclamation points and question marks. Thus, the **<punct>** feature is complex (like the **<agr>** feature), yielding feature equations like **<Punct bal = paren>** or **<Punct term = excl>**. Separating and terminal punctuation marks do not occur adjacent to other members of the same class, but may occasionally occur adjacent to members of the other class, e.g. a question mark on a clause which is separated by a dash from a second clause. Balanced punctuation marks are sometimes adjacent to one another, e.g. quotes immediately inside of parentheses. The **<punct>** feature allows us to control these local interactions.

We also need to control non-local interaction of punctuation marks. Two cases of this are so-called quote alternation, wherein embedded quotation marks must alternate between single and double, and the impossibility of embedding an item containing a colon inside of another item containing a colon. Thus, we have a fourth value for **<punct>**, **<contains colon/dquote/etc. +/->**, which indicates whether or not a constituent contains a particular punctuation mark. This feature is percolated through all auxiliary trees. Things which may not embed are: colons under colons, semi-colons, dashes or commas; semi-colons under semi-

colon or commas. Although it is rare, parentheses may appear inside of parentheses, say with a bibliographic reference inside a parenthesized sentence.

23.1 Appositives, parentheticals and vocatives

These trees handle constructions where additional lexical material is only licensed in conjunction with particular punctuation marks. Since the lexical material is unconstrained (virtually any noun can occur as an appositive), the punctuation marks are anchors and the other nodes are substitution sites. There are cases where the lexical material is restricted, as with parenthetical adverbs like *however*, and in those cases we have the adverb as the anchor and the punctuation marks as substitution sites.

When these constructions can appear inside of clauses (non-peripherally), they must be separated by punctuation marks on both sides. However, when they occur peripherally they have either a preceding or following punctuation mark. We handle this by having both peripheral and non-peripheral trees for the relevant constructions. The alternative is to insert the second (following) punctuation mark in the tokenization process (i.e. insert a comma before the period when an appositive appears on the last NP of a sentence). However, this is very difficult to do accurately.

23.1.1 β nxPUnxPU

The symmetric (non-peripheral) tree for NP appositives, anchored by: comma, dash or parentheses. It is shown in Figure 23.1 anchored by parentheses.

(452) The music here , Russell Smith’s “Tetrameron ” , sounded good . [Brown:cc09]

(453) ...cost 2 million pounds (3 million dollars)

(454) Sen. David Boren (D., Okla.)...

(455) ...some analysts believe the two recent natural disasters – Hurricane Hugo and the San Francisco earthquake – will carry economic ramifications.... [WSJ]

The punctuation marks are the anchors and the appositive NP is substituted. The appositive can be conjoined, but only with a lexical conjunction (not with a comma). Appositives with commas or dashes cannot be pronouns, although they may be conjuncts containing pronouns. When used with parentheses this tree actually presents an alternative rather than an appositive, so a pronoun is possible. Finally, the appositive position is restricted to having nominative or accusative case to block PRO from appearing here.

Appositives can be embedded, as in (456), but do not seem to be able to stack on a single NP. In this they are more like restrictive relatives than appositive relatives, which typically can stack.

(456) ...noted Simon Briscoe, UK economist for Midland Montagu, a unit of Midland Bank PLC.

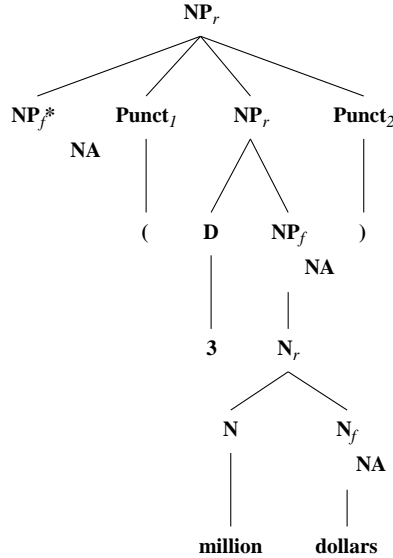


Figure 23.1: The β_{nPxPUxPU} tree, anchored by parentheses

23.1.2 β_{nPxPUxPU}

The symmetric (non-peripheral) tree for N-level NP appositives, is anchored by comma. The modifier is typically an address. It is clear from examples such as (457) that these are attached at N, rather than NP. *Carrier* is not an appositive on *Menlo Park*, as it would be if these were simply stacked appositives. Rather, *Calif.* modifies *Menlo Park*, and that entire complex is compounded with *carrier*, as shown in the correct derivation in Figure 23.2. Because this distinction is less clear when the modifier is peripheral (e.g. ends the sentence), and it would be difficult to distinguish between NP and N attachment, we do not currently allow a peripheral N-level attachment.

(457) An official at Consolidated Freightways Inc., a Menlo Park, Calif, less-than-truckload carrier , said...

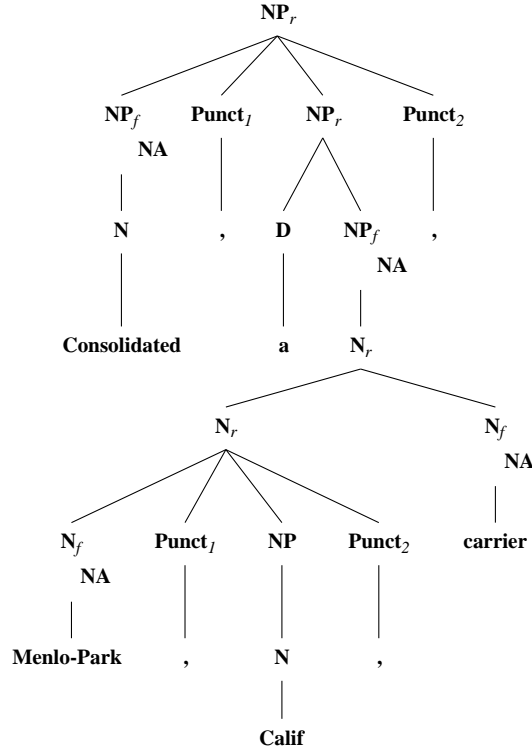
(458) Rep. Ronnie Flippo (D., Ala.), of the delegation, says...

23.1.3 β_{nPxPUx}

This tree, which can be anchored by a comma, dash or colon, handles asymmetric (peripheral) NP appositives and NP colon expansions of NPs. Figure 23.3 shows this tree anchored by a dash and a colon. Like the symmetric appositive tree, β_{nPxPUxpu} , the asymmetric appositive cannot be a pronoun, while the colon expansion can. Thus, this constraint comes from the syntactic entry in both cases rather than being built into the tree.

(459) the bank's 90% shareholder – Petroliam Nasional Bhd. [Brown]

(460) ...said Chris Dillow, senior U.K. economist at Nomura Research Institute .


 Figure 23.2: An N-level modifier, using the β nPUnx tree

(461) ...qualities that are seldom found in one work: Scrupulous scholarship, a fund of personal experience,... [Brown:cc06]

(462) I had eyes for only one person : him .

The colon expansion cannot itself contain a colon, so the foot S has the feature $\text{NP.t} < \text{punctcontainscolon} \geq -$.

23.1.4 β PUpxPUvx

Tree for pre-VP parenthetical PP, anchored by commas or dashes -

(463) John , in a fit of anger , broke the vase

(464) Mary , just within the last year , has totalled two cars

These are clearly not NP modifiers.

Figures 23.4 and 23.5 show this tree alone and as part of the parse for (463).

23.1.5 β puARBpuvx

Parenthetical adverbs - *however*, *though*, etc. Since the class of adverbs is highly restricted, this tree is anchored by the adverb and the punctuation marks substitute. The punctuation marks

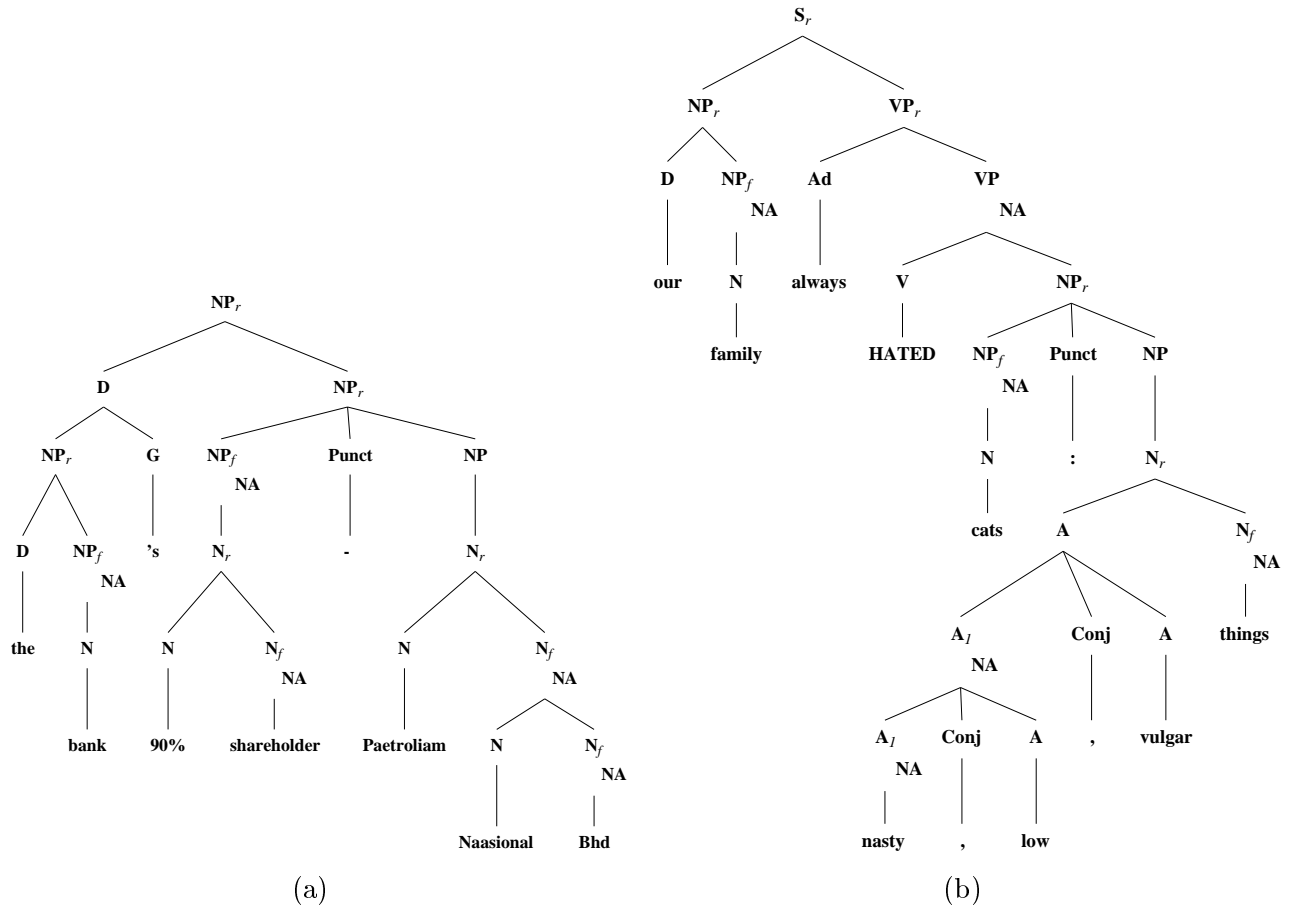


Figure 23.3: The derived trees for an NP with (a) a peripheral, dash-separated appositive and (b) an NP colon expansion (uttered by the Mouse in *Alice's Adventures in Wonderland*)

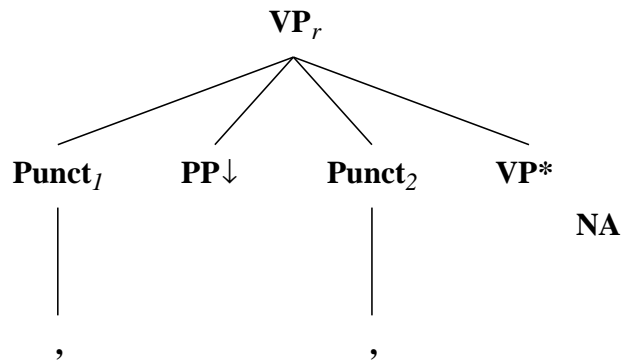
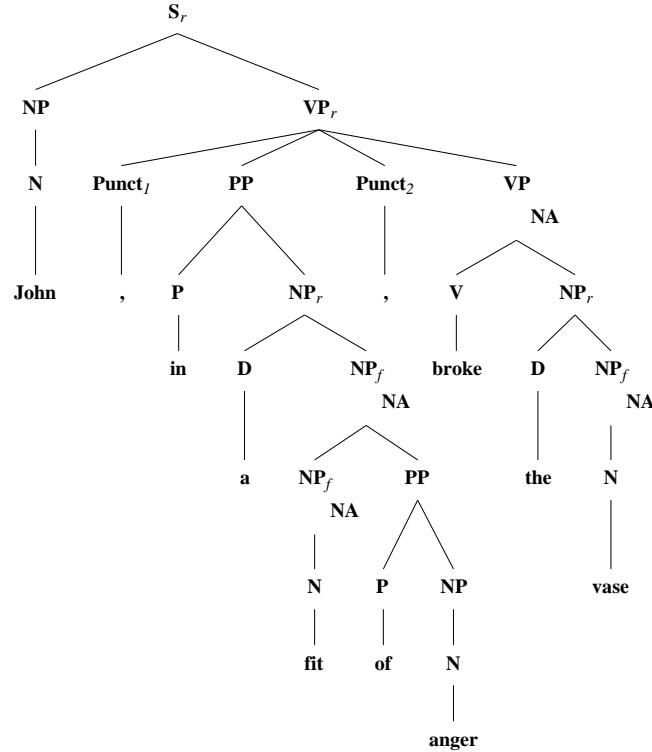


Figure 23.4: The β PUpxPUvx tree, anchored by commas

may be either commas or dashes. Like the parenthetical PP above, these are clearly not NP modifiers.


 Figure 23.5: Tree illustrating the use of $\beta\text{PU}_{\text{pxPUvx}}$

- (465) The new argument over the notification guideline , however , could sour any atmosphere of cooperation that existed . [WSJ]

23.1.6 $\beta\text{sPU}_{\text{nx}}$

Sentence final vocative, anchored by comma:

- (466) You were there , Stanley/my boy .

Also, when anchored by colon, NP expansion on S. These often appear to be extraposed modifiers of some internal NP. The NP must be quite heavy, and is usually a list:

- (467) Of the major expansions in 1960, three were financed under the R. I. Industrial Building Authority's 100mortgage plan: Collyer Wire, Leesona Corporation, and American Tube & Controls.

A simplified version of this sentence is shown in figure 23.6. The NP cannot be a pronoun in either of these cases. Both vocatives and colon expansions are restricted to appear on tensed clauses (indicative or imperative).

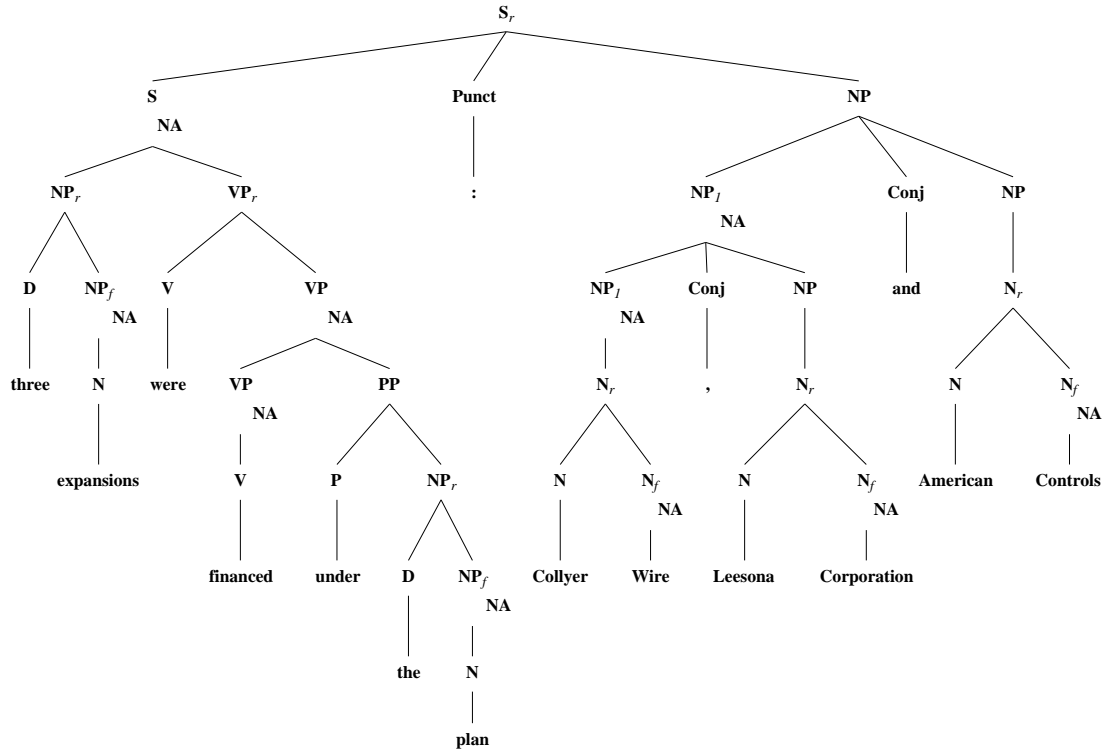


Figure 23.6: A tree illustrating the use of sPUnx for a colon expansion attached at S.

23.1.7 β_{nxPU} s

Tree for sentence initial vocatives, anchored by a comma:

(468) Stanley/my boy , you were there .

The noun phrase may be anything but a pronoun, although it is most commonly a proper noun. The clause adjoined to must be indicative or imperative.

23.2 Bracketing punctuation

23.2.1 Simple bracketing

Trees: $\beta_{\text{PU}}\text{sPU}$, $\beta_{\text{PU}}\text{nxPU}$, $\beta_{\text{PU}}\text{nPU}$, $\beta_{\text{PU}}\text{vxPU}$, $\beta_{\text{PU}}\text{vPU}$, $\beta_{\text{PU}}\text{arbPU}$, $\beta_{\text{PU}}\text{aPU}$, $\beta_{\text{PU}}\text{dPU}$, $\beta_{\text{PU}}\text{pxPU}$, $\beta_{\text{PU}}\text{pPU}$

These trees are selected by parentheses and quotes and can adjoin onto any node type, whether a head or a phrasal constituent. This handles things in parentheses or quotes which are syntactically integrated into the surrounding context. Figure 23.7 shows the $\beta_{\text{PU}}\text{sPU}$ anchored by parentheses, and this tree along with $\beta_{\text{PU}}\text{nxPU}$ in a derived tree.

(469) Dick Carroll and his accordion (which we now refer to as “Freida”) held over at Bahia Cabana where “Sir” Judson Smith brings in his calypso capers Oct. 13 . [Brown:ca31]

- (470) ...noted that the term “teacher-employee” (as opposed to, e.g., “maintenance employee”) was a not inapt description. [Brown:ca35]

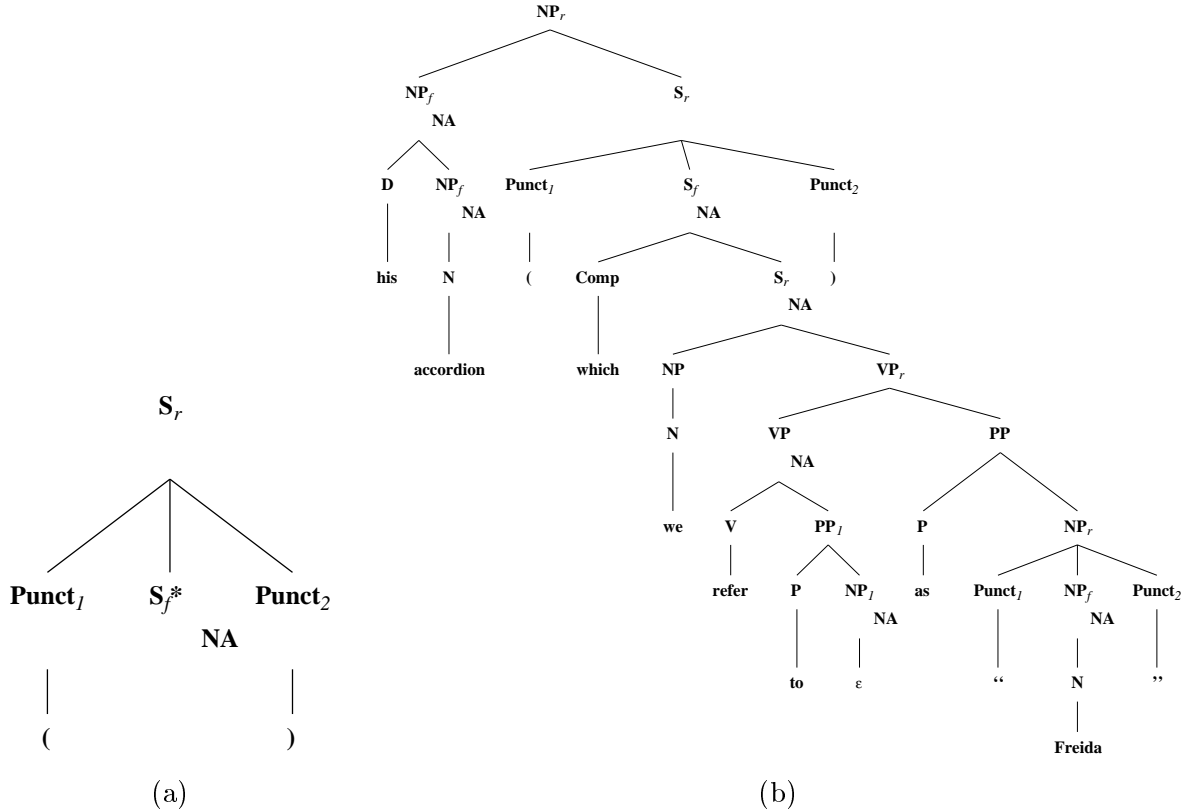


Figure 23.7: β PU_sPU anchored by parentheses, and in a derivation, along with β PU_nxPU

There is a convention in English that quotes embedded in quotes alternate between single and double; in American English the outermost are double quotes, while in British English they are single. The **contains** feature is used to control this alternation. The trees anchored by double quotation marks have the feature **punct contains dquote** = - on the foot node and the feature **punct contains dquote** = + on the root. All adjunction trees are transparent to the **contains** feature, so if any tree below the double quote is itself enclosed in double quotes the derivation will fail. Likewise with the trees anchored by single quotes. The quote trees in effect “toggle” the **contains Xquote** feature. Immediate proximity is handled by the **punct balanced** feature, which allows quotes inside of parentheses, but not vice-versa.

In addition, American English typically places/moves periods (and commas) inside of quotation marks when they would logically occur outside, as in example 471. The comma in the first part of the quote is not part of the quote, but rather part of the parenthetical quoting clause. However, by convention it is shifted inside the quote, as is the final period. British English does not do this. We assume here that the input has already been tokenized into the “British” format.

- (471) “You can’t do this to us ,” Diane screamed . “We are Americans.”

The β PU s PU can handle quotation marks around multiple sentences, since the sPUs tree allows us to join two sentences with a period, exclamation point or question mark. Currently, however, we cannot handle the style where only an open quote appears at the beginning of a paragraph when the quotation extends over multiple paragraphs. We could allow a lone open quote to select the β PU s tree, if this is deemed desirable.

Also, the β PU s PU is selected by a pair of commas to handle non-peripheral appositive relative clauses, such as in example (472). Restrictive and appositive relative clauses are not syntactically differentiated in the XTAG grammar (cf. Chapter 14).

(472) This news , announced by Jerome Toobin , the orchestra’s administrative director , brought applause ... [Brown:cc09]

The trees discussed in this section will only allow balanced punctuation marks to adjoin to constituents. We will not get them around non-constituents, as in (473).

(473) Mary asked him to leave (and he left)

23.2.2 β sPU s PU

This tree allows a parenthesized clause to adjoin onto a non-parenthesized clause.

(474) Innumerable motels from Tucson to New York boast swimming pools (“ swim at your own risk ” is the hospitable sign poised at the brink of most pools) . [Brown:ca17]

23.3 Punctuation trees containing no lexical material

23.3.1 α PU

This is the elementary tree for substitution of punctuation marks. This tree is used in the quoted speech trees, where including the punctuation mark as an anchor along with the verb of saying would require a new entry for every tree selecting the relevant tree families. It is also used in the tree for parenthetical adverbs (β puARBpuvx), and for S-adjoined PPs and adverbs (β spuARB and β spuPnx).

23.3.2 β PU s

Anchored by comma: allows comma-separated clause initial adjuncts, (475-476).

(475) Here , as in “Journal” , Mr. Louis has given himself the lion’s share of the dancing... [Brown:cc09]

(476) Choreographed by Mr. Nagrin, the work filled the second half of a program

To keep this tree from appearing on root S s (i.e. , *sentence*), we have a root constraint that **<punct struct = nil>** (similar to the requirement that root S s be tensed, i.e. **<mode = ind/imp>**). The **<punct struct> = nil** feature on the foot blocks stacking of multiple punctuation marks. This feature is shown in the tree in Figure 23.8.

This tree can be also used by adjuncts on embedded clauses:

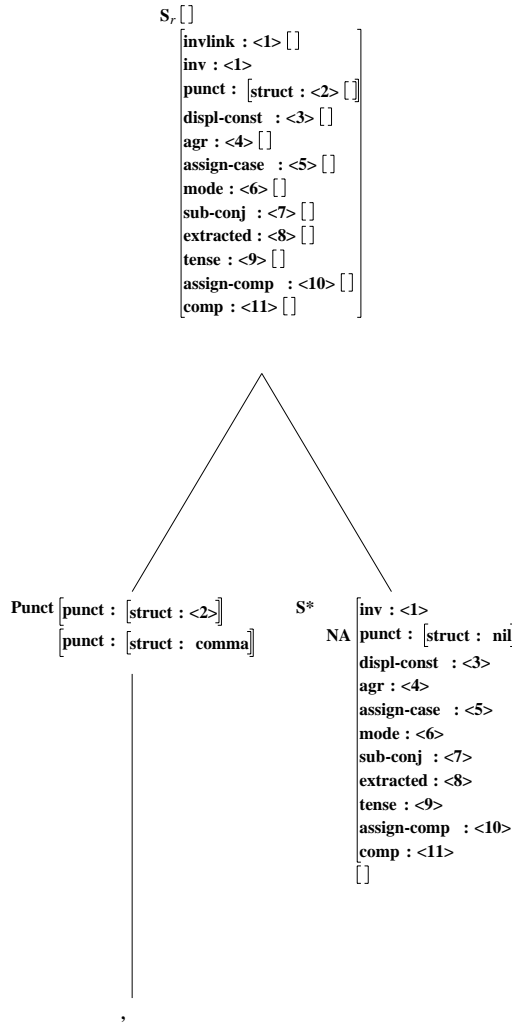


Figure 23.8: β PU, with features displayed

- (477) One might expect that in a poetic career of seventy-odd years, some changes in style and method would have occurred, some development taken place. [Brown:cj65]

These adjuncts sometimes have commas on both sides of the adjunct, or, like (477), only have them at the end of the adjunct.

Finally, this tree is also used for peripheral appositive relative clauses.

- (478) Interest may remain limited into tomorrow's U.K. trade figures, which the market will be watching closely to see if there is any improvement after disappointing numbers in the previous two months.

23.3.3 β sPUs

This tree handles clausal “coordination” with comma, dash, colon, semi-colon or any of the terminal punctuation marks. The first clause must be either indicative or imperative. The second may also be infinitival with the separating punctuation marks, but must be indicative or imperative with the terminal marks; with a comma, it may only be indicative. The two clauses need not share the same mode. NB: Allowing the terminal punctuation marks to anchor this tree allows us to parse sequences of multiple sentences. This is not the usual mode of parsing; if it were, this sort of sequencing might be better handled by a higher level of processing.

(479) For critics , Hardy has had no poetic periods – one does not speak of early Hardy or late Hardy , or of the London or Max Gate period....

(480) Then there was exercise , boating and hiking , which was not only good for you but also made you more virile : the thought of strenuous activity left him exhausted.

This construction is one of the few where two non-bracketing punctuation marks can be adjacent. It is possible (if rare) for the first clause to end with a question mark or exclamation point, when the two clauses are conjoined with a semi-colon, colon or dash. Features on the foot node, as shown in Figure 23.9, control this interaction.

Complementizers are not permitted on either conjunct. Subordinating conjunctions sometimes appear on the right conjunct, but seem to be impossible on the left:

(481) Killpath would just have to go out and drag Gun back by the heels once an hour ; because he’d be damned if he was going to be a mid-watch pencil-pusher . [Brown:cl17]

(482) The best rule of thumb for detecting corked wine (provided the eye has not already spotted it) is to smell the wet end of the cork after pulling it : if it smells of wine , the bottle is probably all right ; if it smells of cork , one has grounds for suspicion. [Brown:cf27]

23.3.4 β sPU

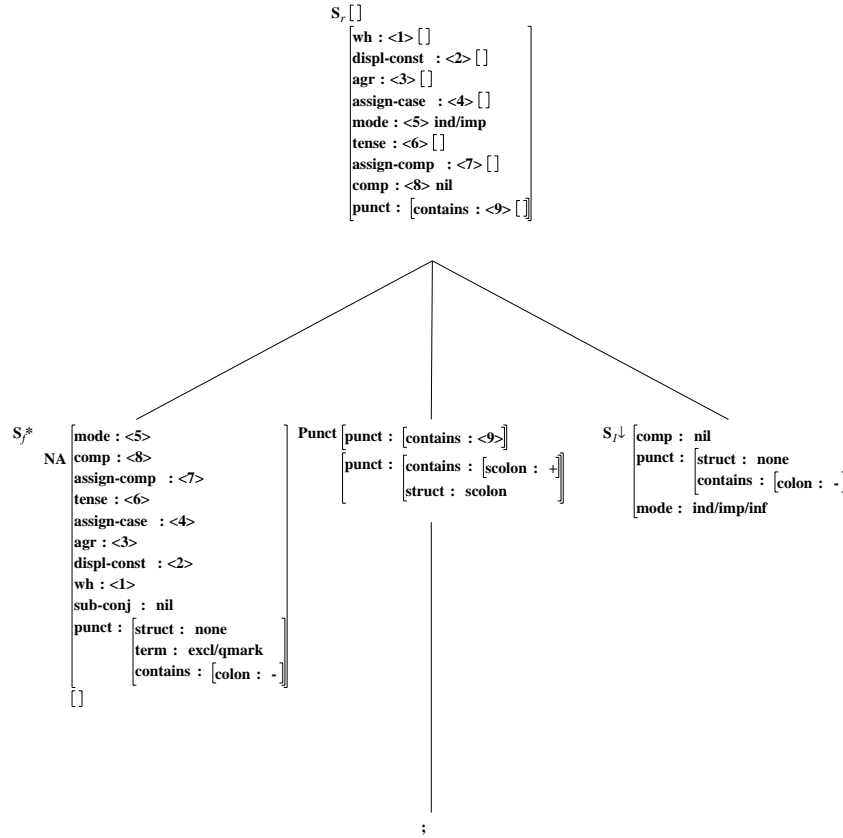
This tree handles the sentence final punctuation marks when selected by a question mark, exclamation point or period. One could also require a final punctuation mark for all clauses, but such an approach would not allow non-periods to occur internally, for instance before a semi-colon or dash as noted above in Section 23.3.3. This tree currently only adjoins to indicative or imperative (root) clauses.

(483) He left !

(484) Get lost .

(485) Get lost ?

The feature **punct bal= nil** on the foot node ensures that this tree only adjoins inside of parentheses or quotes completely enclosing a sentence (486), but does not restrict it from adjoining to clause which ends with balanced punctuation if only the end of the clause is contained in the parentheses or quotes (487).


 Figure 23.9: β sPUs, with features displayed

(486) (John then left .)

(487) (John then left) .

(488) Mary asked him to leave (immediately) .

This tree is also selected by the colon to handle a colon expansion after adjunct clause –

(489) Expressed differently : if the price for becoming a faithful follower... [Brown:cd02]

(490) Expressing it differently : if the price for becoming a faithful follower...

(491) To express it differently : if the price for becoming a faithful follower... [Brown:cd02]

This tree is only used after adjunct (untensed) clauses, which adjoin to the tensed clause using the adjunct clause trees (cf Section 15); the **mode** of the complete clause is that of the matrix rather than the adjunct. Indicative or imperative (i.e. root) clauses separated by a colon use the β sPUs tree (Section 23.3.3).

23.3.5 β_v PU

This tree is anchored by a colon or a dash, and occurs between a verb and its complement. These typically are lists.

- (492) Printed material Available , on request , from U.S. Department of Agriculture , Washington 25 , D.C. , are : Cooperative Farm Credit Can Assist..... [Brown:ch01]

23.3.6 β_p PU

This tree is anchored by a colon or a dash, and occurs between a preposition and its complement. It typically occurs with a sequence of complements. As with the tree above, this typically occurs with a conjoined complement.

- (493) ...and utilization such as : (A) the protection of forage...

- (494) ...can be represented as : Af.

23.4 Other trees

23.4.1 β_{spu} ARB

In general, we attach post-clausal modifiers at the VP node, as you typically get scope ambiguity effects with negation (*John didn't leave today* – did he leave or not?). However, with post-sentential, comma-separated adverbs, there is no ambiguity - in *John didn't leave, today* he definitely did not leave. Since this tree is only selected by a subset of the adverbs (namely, those which can appear pre-sententially, without a punctuation mark), it is anchored by the adverb.

- (495) The names of some of these products don't suggest the risk involved in buying them , either . [WSJ]

23.4.2 β_{spu} Pnx

Clause-final PP separated by a comma. Like the adverbs described above, these differ from VP adjoined PPs in taking widest scope.

- (496) ...gold for current delivery settled at \$367.30 an ounce , up 20 cents .

- (497) It increases employee commitment to the company , with all that means for efficiency and quality control .

23.4.3 β_{nx} PUa

Anchored by colon or dash, allows for post-modification of NPs by adjectives.

- (498) Make no mistake , this Gorky Studio drama is a respectable import – aptly grave , carefully written , performed and directed .

CHAPTER 23. PUNCTUATION MARKS

Part V

Appendices

Appendix A

Future Work

A.1 Adjective ordering

At this point, the treatment of adjectives in the XTAG English grammar does not include selectional or ordering restrictions.¹ Consequently, any adjective can adjoin onto any noun and on top of any other adjective already modifying a noun. All of the modified noun phrases shown in (499)-(502) currently parse.

(499) big green bugs

(500) big green ideas

(501) colorless green ideas

(502) *green big ideas

While (500)-(502) are all semantically anomalous, (502) also suffers from an ordering problem that makes it seem ungrammatical as well. Since the XTAG grammar focuses on syntactic constructions, it should accept (499)-(501) but not (502). Both the auxiliary and determiner ordering systems are structured on the idea that certain types of lexical items (specified by features) can adjoin onto some types of lexical items, but not others. We believe that an analysis of adjectival ordering would follow the same type of mechanism.

A.2 More work on Determiners

In addition to the analysis described in Chapter 18, there remains work to be done to complete the analysis of determiner constructions in English.² Although constructions such as determiner coordination are easily handled if overgeneration is allowed, blocking sequences such as *one and some* while allowing sequences such as *five or ten* still remains to be worked out. There are still a handful of determiners that are not currently handled by our system. We do not have

¹This section is a repeat of information found in section 19.1.

²This section is from [Hockey and Mateyak, 1998].

an analysis to handle *most*, *such*, *certain*, *other* and *own*³. In addition, there is a set of lexical items that we consider adjectives (*enough*, *less*, *more* and *much*) that have the property that they cannot cooccur with determiners. We feel that a complete analysis of determiners should be able to account for this phenomenon, as well.

A.3 *-ing* adjectives

An analysis has already been provided for past participial (*-ed*) adjectives (as in sentence (503)), which are restricted to the Transitive Verb family.⁴ A similar analysis needs to take place for the present participle (*-ing*) used as a pre-nominal modifier. This type of adjective, however, does not seem to be as restricted as the *-ed* adjectives, since verbs in other tree families seem to exhibit this alternation as well (e.g. sentences (504) and (505)).

(503) The murdered man was a doctoral student at UPenn .

(504) The man died .

(505) The dying man pleaded for his life .

A.4 Verb selectional restrictions

Although we explicitly do not want to model semantics in the XTAG grammar, there is some work along the syntax/semantics interface that would help reduce syntactic ambiguity and thus decrease the number of semantically anomalous parses. In particular, verb selectional restrictions, particularly for PP arguments and adjuncts, would be quite useful. With the exception of the required *to* in the Ditransitive with PP Shift tree family (Tnx0Vnx1Pnx2), any preposition is allowed in the tree families that have prepositions as their arguments. In addition, there are no restrictions as to which prepositions are allowed to adjoin onto a given verb. The sentences in (506)-(508) are all currently accepted by the XTAG grammar. Their violations are stronger than would be expected from purely semantic violations, however, and the presence of verb selectional restrictions on PP's would keep these sentences from being accepted.

(506) #survivors walked of the street .

(507) #The man about the earthquake survived .

(508) #The president arranged on a meeting .

³The behavior of *own* is sufficiently unlike other determiners that it most likely needs a tree of its own, adjoining onto the right-hand side of genitive determiners.

⁴This analysis may need to be extended to the Transitive Verb particle family as well.

A.5 Thematic Roles

Elementary trees in TAGs capture several notions of locality, with the most primary of these being locality of θ -role assignment. Each elementary tree has associated with it the θ -roles assigned by the anchor of that elementary tree. In the current XTAG system, while the notion of locality of θ -role assignment within an elementary tree has been implicit, the θ -roles assigned by a head have not been explicitly represented in the elementary tree. Incorporating θ -role information will make the elementary trees more informative and will enable efficient pruning of spurious derivations when embedded into a specific context. In the case of a Synchronous TAG, θ -roles can also be used to automatically establish links between two elementary trees, one in the object language and one in the target language.

Appendix B

Metarules

B.1 Introduction

XTAG has now a collection of functions accessible from the user interface that helps the user in the construction and maintenance of a tag tree-grammar. This subsystem is based on the idea of metarules ([Becker, 1993]). Here our primary purpose is to describe the facilities implemented under this metarule-based subsystem. For a discussion of the metarules as a method for compact representation of the Lexicon see [Becker, 1993] and [Srinivas *et al.*, 1994].

The basic idea of using metarules is to take profit of the similarities of the relations involving related pairs of XTAG elementary trees. For example, in the English grammar described in this technical report, comparing the XTAG trees for the basic form and the wh-subject moved form, the relation between this two trees for transitive verbs ($\alpha nx_0 V nx_1$, $\alpha W_0 nx_0 V nx_1$) is similar to the relation for the intransitive verbs ($\alpha nx_0 V$, $\alpha W_0 nx_0 V$) and also to the relation for the ditransitives ($\alpha nx_0 V nx_1 nx_2$, $\alpha W_0 nx_0 V nx_1 nx_2$). Hence, instead of generating by hand the six trees mentioned above, a more natural and robust way would be generating by hand only the basic trees for the intransitive, transitive and ditransitive cases, and letting the wh-subject moved trees to be automatically generated by the application of a unique transformation rule that would account exactly for the identical relation involved in each of the three pairs above.

Notice that the degree of generalization can be much higher than it might be thought in principle from the above paragraph. For example, once a rule for passivization is applied to the tree different basic trees above, the wh-subject moved rule could be again applied to generate the wh-moved subject versions for the passive form. Depending on the degree of regularity that one can find in the grammar being built, the reduction in the number of original trees can be exponential.

We still make here a point that the reduction of effort in grammar construction is not the only advantage of the approach. Robustness, reliability and maintainability of the grammar achieved by the use of metarules are equally or even more important.

In the next section we define a metarule in XTAG. Section 3 gives some linguistically motivated examples of metarule for the English grammar described in this technical report and their application. Section 4 describes the access through the user interface.

B.2 The definition of a metarule in XTAG

A metarule specifies a rule for transforming grammar rules into grammar rules. In XTAG the grammar rules are lexicalized trees. Hence an XTAG metarule **mr** is a pair (**lhs**, **rhs**) of XTAG trees, where:

- **lhs**, the *left-hand side* of the metarule, is a pattern tree, i.e., it is intended to present a specific pattern of tree to look for in the trees submitted to the application of the metarule.
- When a metarule **mr** is applied to an input tree **inp**, the first step is to verify if the input tree matches the pattern specified by the **lhs**. If there is no match, the application *fails*.
- **rhs**, the *right-hand side* of the metarule, specifies (together with **lhs**) the transformation that will be done in **inp**, in case of successful matching, thus generating the output tree of the metarule application¹.

B.2.1 Node names, variable instantiation, and matches

We will use the terms (**lhs**, **rhs** and **inp**) as introduced above to refer to the parts of a generic metarule being applied to an input tree.

The nodes at **lhs** can take three different forms: a constant node, a typed variable node, and a non-typed variable node. The naming conventions for these different classes of nodes is given below.

- **Constant Node:** Its name must not initiate by a question mark ('?' character). They are like we expect for names to be in normal XTAG trees; for instance, **inp** is expected to have only constant nodes. Some examples of constant nodes are NP , V , NP_0 , NP_1 , S_r . We will call the two parts that compose such names the *stem* and the *subscript*. In the examples above NP , V and S are stems and 0, 1, r are subscripts. Notice that the subscript part can also be empty as in two of the above examples.
- **Non-Typed Variable Node:** Its name initiates by a question mark ('?'), followed by a sequence of digits (i.e. a number) which uniquely identifies the variable. Examples: ?1, ?3, ?3452². We assume that there is no stem and no subscript in this names, i.e., '?' is just a meta-character to introduce a variable, and the number is the variable identifier.
- **Typed Variable Node:** Its name initiates by a question mark ('?') followed by a sequence of digits, but is additionally followed by a *type specifiers definition*. A *type specifiers definition* is a sequence of one or more *type specifier* separated by a slash ('/'). A *type specifier* has the same form of a regular XTAG node name (like the constant nodes), except that the subscript can be also a question mark. Examples of typed variables are: ?1 VP (a single type specifier with stem VP and no subscript), ?3 NP_1/PP (two type specifiers, NP_1 and PP), ?1 $NP?$ (one type specifier, $NP?$ with undetermined subscript).

¹actually more than one output tree can be generated from the successful application of a rule to an input tree, as will be seen soon

²Notice however that having the sole purpose of distinguishing between variables, a number like the one in the last example is not very likely to occur, and a metarule with more than three thousand variables can give you a place in the Guinness TagBook of Records.

We'll see ahead that each type specifier represents an alternative for matching, and the presence of '?' in subscript position of a type specifier means that matching will only check for the stem ³.

During the process of matching, variables are associated (we use the term *instantiated*) with 'tree material'. According to its class a variable can be instantiated with different kinds of tree material:

- A typed variable will be instantiated with exactly one node of the input tree, which is in accordance to one of its type specifiers (The full rule is in the following subsection).
- A non-typed variable will be instantiated with a range of subtrees. These subtrees will be taken from one of the nodes of the input tree **inp**. Hence, there will a node n in **inp**, with subtrees $n.t_1, n.t_2, \dots, n.t_k$, in this order, where the variable will be instantiated with some subsequence of these subtrees (e.g., $n.t_2, n.t_3, n.t_4$). Note however, that some of these subtrees, may be incomplete, i.e., they may not go all the way to the bottom leaves. Entire subtrees may be removed. Actually for each child of the non-typed variable node, one subtree that matches this child subtree will be removed from some of the $n.t_i$ (maybe an entire $n.t_i$), leaving in place a mark for inserting material during the substitution of occurrences at **rhs**.

Notice still that the variable can be instantiated with a single tree and even with no tree.

We define a *match* to be a complete instantiation of all variables appearing in the metarule. In the process of matching, there may be several possible ways of instantiating the set of variables of the metarule, i.e., several possible matches. This is due to the presence of non-typed variables.

Now, we are ready to define what we mean by a successful matching. The process of matching is *successful* if the number of possible matches is greater then 0. When there is no possible match the process is said to *fail*. In addition to return success or failure, the process also return the set of all possible *matches*, which will be used for generating the output.

B.2.2 Structural Matching

The process of matching **lhs** and **inp** can be seen as a recursive procedure for matching trees, starting at their roots and proceeding in a top-down style along with their subtrees. In the explanation of this process that follows we have used the term **lhs** not only to refer to the whole tree that contains the pattern but to any of its subtrees that is being considered in a given recursive step. The same applies to **inp**. By now we ignore feature equations, which will be accounted for in the next subsection.

The process described below returns at the end the set of matches (where an empty set means the same as failure). We first give one auxiliary definition, of valid Mapping, and one recursive function Match, that matches lists of trees instead of trees, and then define the process of matching two trees as a special case of call to Match.

³This is different from not having a subscript which is interpreted as checking that the input tree have no subscript for matching

Given a list $list_{lhs} = [lhs_1, lhs_2, \dots, lhs_l]$ of nodes of **lhs** and a list $list_{inp} = [inp_1, inp_2, \dots, inp_i]$ of nodes of **inp**, we define a *mapping* from $list_{lhs}$ to $list_{inp}$ to be a function *Mapping*, that for each element of $list_{lhs}$ assigns a list of elements of $list_{inp}$, defined by the following condition:

$$concatenation (Mapping(lhs_1), Mapping(lhs_2), \dots, Mapping(lhs_l)) = list_{inp}$$

, i.e., the elements of $list_{inp}$ are split into sublists and assigned in order of appearance in the list to the elements of $list_{lhs}$.

We say that a mapping is a *valid mapping* if for all j , $1 \leq j \leq l$ (where l is the length of $list_{lhs}$), the following restrictions apply:

1. if lhs_j is a constant node, then $Mapping(lhs_j)$ must have a single element, say, $rhs_{g(j)}$, and the two nodes must have the same name and agree on the markers (foot, substitution, head and NA), i.e., if lhs_j is NA, then $rhs_{g(j)}$ must be NA, if lhs_j has no markers, then $rhs_{g(j)}$ must have no markers, etc.
2. if lhs_j is a type variable node, then $Mapping(lhs_j)$ must have a single element, say, $rhs_{g(j)}$, and $rhs_{g(j)}$ must be *marker-compatible* and *type-compatible* with lhs_j .
 $rhs_{g(j)}$ is *marker-compatible* with lhs_j if any marker (foot, substitution, head and NA) present in lhs_j is also present in $rhs_{g(j)}$ ⁴.
 $rhs_{g(j)}$ is *type-compatible* with lhs_j if there is at least one of the alternative type specifiers for the typed variable that satisfies the conditions below.
 - $rhs_{g(j)}$ has the stem defined in the type specifier.
 - if the type specifier doesn't have subscript, then $rhs_{g(j)}$ must have no subscript.
 - if the type specifier has a subscript different of '?', then $rhs_{g(j)}$ must have the same subscript as in the type specifier⁵.
3. if lhs_j is a non-typed variable node, then there's actually no requirement: $Mapping(lhs_j)$ may have any length and even be empty.

The following algorithm, Match, takes as input a list of nodes of **lhs** and a list of nodes of **inp**, and returns the set of possible matches generated in the attempt of match this two lists. If the result is an empty set, this means that the matching failed.

Function Match ($list_{lhs}, list_{rhs}$)

Let *MAPPINGS* be the list of all valid mappings from $list_{lhs}$ to $list_{rhs}$

Make *MATCHES* = \emptyset

For each mapping $Mapping \in MAPPINGS$ do:

Make *Matches* = $\{\emptyset\}$

For each j , $1 \leq j \leq l$, where $l = length(list_{lhs})$, do:

if lhs_j is a constant node, then

⁴Notice that, unlike the case for the constant node, the inverse is not required, i.e., if lhs_j has no marker, $rhs_{g(j)}$ is still allowed to have some.

⁵If the type specifier has a '?' subscript, there is no restriction, and that is exactly its function: to allow for the matching to be independent of the subscript

let $children_{lhs}$ be the list of children of lhs_j
 $lhr_{g(j)}$ be the single element in $Mapping(lhs_j)$
 $children_{rhs}$ be the list of children of $lhr_{g(j)}$
 Make $Matches = \{m \cup m_j \mid m \in Matches$
 $\quad \text{and } m_j \in Match(children_{lhs}, children_{rhs})\}$
 if lhs_j is a typed variable node, then
 let $children_{lhs}$ be the list of children of lhs_j
 $lhr_{g(j)}$ be the single element in $Mapping(lhs_j)$
 $children_{rhs}$ be the list of children of $lhr_{g(j)}$
 Make $Matches = \{(lhs_j, lhr_{g(j)})\} \cup m \cup m_j \mid m \in Matches$
 $\quad \text{and } m_j \in Match(children_{lhs}, children_{rhs})\}$
 if lhs_j is a non-typed variable node, then
 let $children_{lhs}$ be the list of children of lhs_j
 sl be the number of nodes in $children_{lhs}$
 $DESC_s$ be the set of s-size lists given by:
 $DESC_s = \{[dr_1, dr_2, \dots, dr_s] \mid$
 $\quad \text{for every } 1 \leq k \leq s, dr_k \text{ is a descendant}$
 $\quad \text{of some node in } Mapping(lhs_j)^6$
 $\quad \text{for every } 1 < k \leq s, dr_k \text{ is to the right of } dr_{k-1}^7.$
 For every list $Desc = [dr_1, dr_2, \dots, dr_s] \in DESC_s$ do:
 Let Tree-Material be the list of subtrees dominated
 by the nodes in $Mapping(lhs_j)$, but, with the
 subtrees dominated by the nodes in $DESC_s$
 cut out from these trees
 Make $Matches = \{(lhs_j, Tree - Struct)\} \cup m \cup m_j \mid$
 $\quad m \in Matches \text{ and } m_j \in Match(children_{lhs}, Desc)\}$
 Make $MATCHES = MATCHES \cup Matches$
 Return $MATCHES$

Finally we can define the process of structurally matching **lhs** to **inp** as the evaluation of $Match([root(\mathbf{lhs})], [root(\mathbf{inp})])$. If the result is an empty set, the matching failed, otherwise the resulting set is the set of possible matches that will be used for generating the new trees (after being pruned by the feature equation matching).

B.2.3 Output Generation

Although nothing has yet been said about the feature equations, which is the subject of the next subsection, we assume that only matches that meet the additional constraints imposed by feature equations are considered for output. If no structural match survives feature equations checking, that matching has failed.

If the process of matching **lhs** to **inp** fails, there are two alternative behaviors according to the value of a parameter⁸. If the parameter is set to false, which is the *default* value, no output is generated. On the other hand, if it is set to true, then the own **inp** tree is copied to the output.

⁸the parameter is accessible at the Lisp interface by the name *XTAG::*metarules-copy-unmatched-trees**

If the process of matching succeeds, as many trees will be generated in the output as the number of possible matches obtained in the process. For a given match, the output tree is generated by substituting in the **rhs** tree of the metarule the occurrences of variables by the material to which they have been instantiated in the match. The case of the typed-variable is simple. The name of the variable is just substituted by the name of the node to which it has been instantiated from **inp**. A very important detail is that the marker (foot, substitution, head, NA, or none) at the output tree node comes from what is specified in the **rhs** node, which can be different of the marker at the variable node in **inp** and of the associated node from **inp**.

The case of the non-typed variable, not surprisingly, is not so simple. In the output tree, this node will be substituted by the subtree list that was associated to this node, in the same other, attaching to the parent of this non-typed variable node. But remember, that some subtrees may have been removed from some of the trees in this list, maybe entire elements of this list, due to the effect of the children of the metavariable in **lhs**. It is a requirement that any occurrence of a non-typed variable node at the **rhs** tree has exactly the same number of children than the unique occurrence of this non-typed variable node in **lhs**. Hence, when generating the output tree, the subtrees at **rhs** will be inserted exactly at the points where subtrees were removed during matching, in a positional, one to one correspondance.

For feature equations in the output trees see the next subsection. The comments at the output are the comments at the **lhs** tree of the metarule followed by the comments at **inp**, both parts introduced by appropriate headers, allowing the user to have a complete history of each tree.

B.2.4 Feature Matching

In the previous subsections we have considered only the aspects of a metarule involving the structural part of the XTAG trees. In a feature based grammar as XTAG is, accounting for features is essential. A metarule is not really worth if it doesn't account for the proper change of feature equations⁹ from the input to the output tree. The aspects that have to be considered here are:

- Which feature equations should be required to be present in **inp** in order for the match to succeed.
- Which feature equations should be generated in the output tree as a function of the feature equations in the input tree.

Based on the possible combinations of these requirements we partition the feature equations into the following five classes¹⁰:

- *Require & Retain*: Feature equations in this class are required to be in **inp** in order for matching to succeed. Upon matching, these equations will be copied to the output tree.

⁹Notice that what is really important is not the features themselves, but the feature equations that relate the feature values of nodes of the same tree

¹⁰This classification is really a partition, i.e., no equation may be conceptually in more than one class at the same time.

To achieve this behaviour, the equation must be placed in the **lhs** tree of the metarule preceded by a plus character (e.g. $+V.t :< trans > = +$)¹¹

- *Require & Don't Copy*: The equation is required to be in **inp** for matching, but should not be copied to the output tree. Those equations must be in **lhs** preceded by minus character (e.g. $-NP_1 :< case > = acc$).
- *Optional & Don't Copy*: The equation is not required for matching, but we have to make sure not to copy it to the output tree set of equations, regardless of it being present or not in **inp**. Those equations must be in **lhs** in raw form, i.e. neither preceded by a plus nor minus character (e.g. $S_r.b :< perfect > = VP.t :< perfect >$).
- *Optional & Retain*: The equation is not required for matching but, in case it is found in **inp** it must be copied to the output tree. This is the *default* case, and hence these equations should not be present in the metarule specification.
- *Add*: The equation is not required for matching but we want it to be put in the output tree anyway. These equations are placed in raw form in the **rhs** (notice in this case it is the right hand side).

Typed variables can be used in feature equations in both **lhs** and **rhs**. They are intended to represent the nodes of the input tree to which they have been instantiated. For each resulting match from the structural matching process the following is done:

- The (typed) variables in the equations at **lhs** and **rhs** are substituted by the names of the nodes they have been instantiated to.
- The requirements concerning feature equations are checked, according to the above rules.
- If the match survives feature equation checking, the proper output tree is generated, according to Section B.2.3 and to the rules described above for the feature equations.

Finally, a new kind of metavariable, which is not used at the nodes, can be introduced in the feature equations part. They have the same form of the non-typed variables, i.e. quotation mark, followed by a number, and are used in the place of feature values and feature names. Hence, if the equation $?NP_?.b :<?2 > = ?3$ appears in **lhs**, this means, that all feature equations of **inp** that match a bottom attribute of some *NP* to any feature value (but not to a feature path) will not be copied to the output.

B.3 Examples

Figure B.1 shows a metarule for wh-movement of the subject. Among the trees to which it have been applied are the basic trees of intransitive, transitive and ditransitive families (including prepositional complements), passive trees of the same families, and ergative.

¹¹Commutativity of equations is accounted for in the system. Hence an equation $x = y$ can also be specified as $y = x$. Associativity is not accounted for and its need by an user is viewed as indicating misspecification at the input trees.

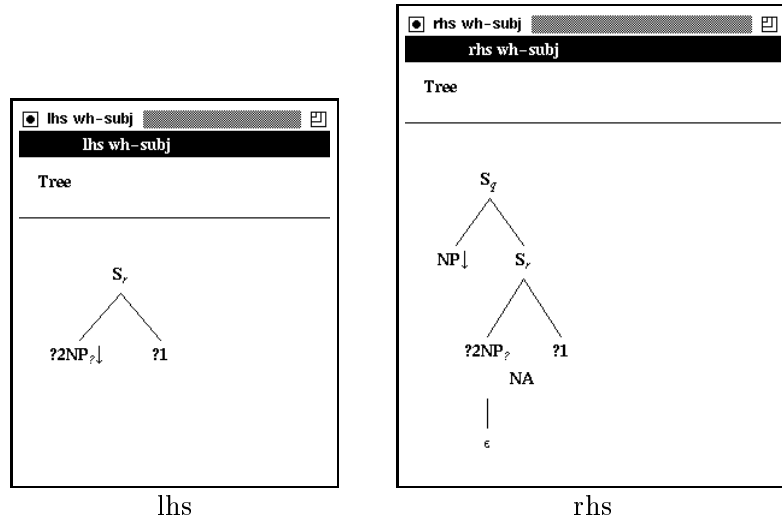


Figure B.1: Metarule for wh-movement of subject

Figure B.2 shows a metarule for wh-movement of an NP in object position. Among the trees to which it have been applied are the basic and passive trees of transitive and ditransitive families.

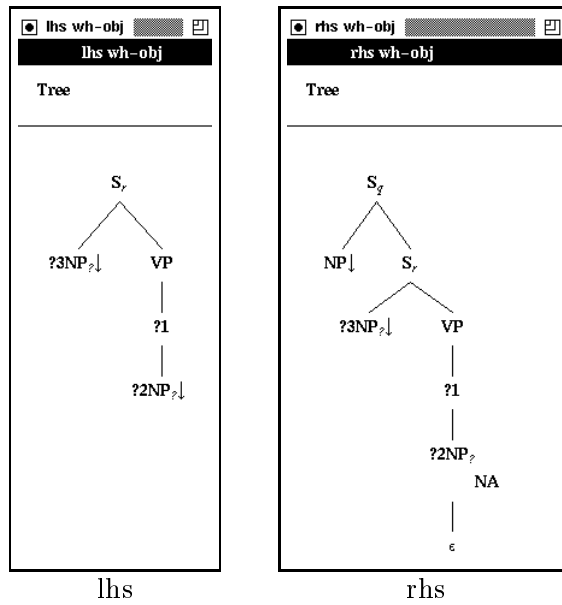


Figure B.2: Metarule for wh-movement of object

Figure B.3 shows a metarule for general wh-movement of an NP. It can be applied to

generate trees with either subject or object NP moved. We show in Figure B.4, the basic tree for the family Tnx0Vnx1Pnx2 and the tree wh-trees generated by the application of the rule.

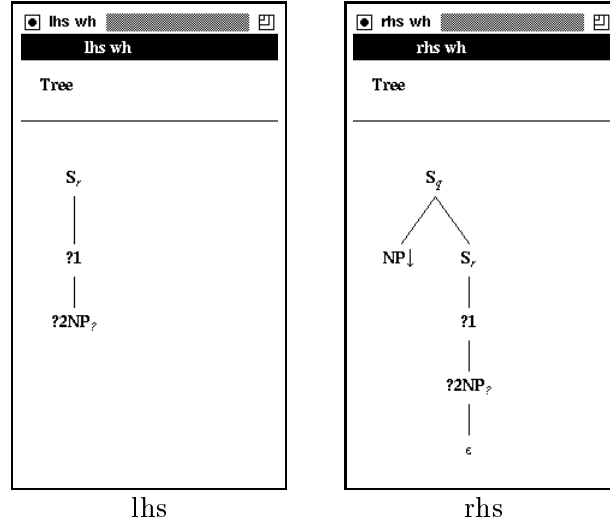


Figure B.3: Metarule for general wh movement of an NP

B.4 The Access to the Metarules through the XTAG Interface

We first describe the access to the metarules subsystem using buffers with single metarule applications. Then we proceed by describing the application of multiple metarules in what we call the parallel, sequential, and cumulative modes to input tree files.

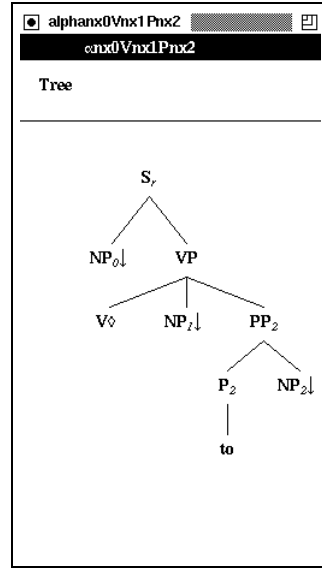
We have defined conceptually a metarule as an ordered pair of trees. In the implementation of the metarule subsystem it works the same: a metarule is a buffer with two trees. The name of the metarule is the name of the buffer. The first tree that appear in the main window under the metarule buffer is the *left hand side*, the next appearing below is the *right hand side*¹². The positional approach allows us to have naming freedom: the tree names are irrelevant¹³. Since we can save buffers into text files, we can talk also about metarule files.

The available options for applying a metarule which is in a buffer are:

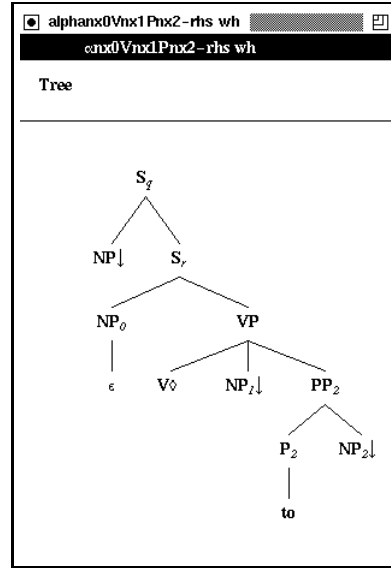
- For applying it to a single input tree, click in the name of the tree in the main window, and choose the option *apply metarule to tree*. You will be prompted for the name of the metarule to apply to the tree which should be, as we mentioned before, the name of the buffer that contains the metarule trees. The output trees will be generated at the end of the buffer that contains the input tree. The names of the trees depend of a LISP parameter **metarules-change-name**. If the value of the parameter is **false** — the *default*

¹²Although a buffer is intended to implement the concept of a set (not a sequence) of trees we take profit of the actual organization of the system to realize the concept of (ordered) tree pair in the implementation.

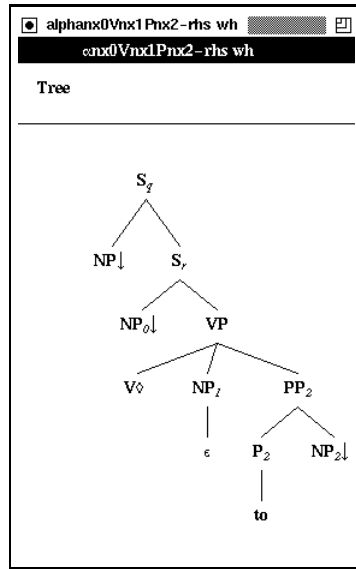
¹³so that even if we want to have mnemonic names resembling their distinct character - left or right hand side, - we have some naming flexibility to call them e.g. *lhs23* or *lhs-passive*, ...



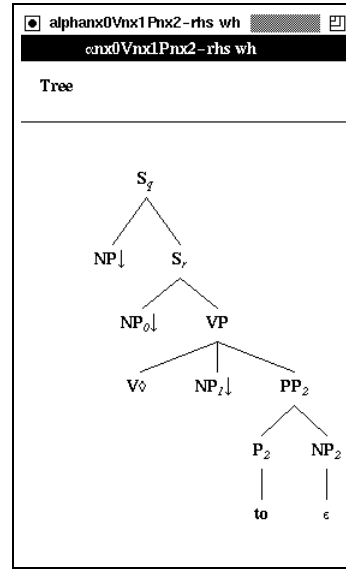
Tnx0Vnx1Pnx2



subject moved



NP object moved



NP object moved from PP

Figure B.4: Application of wh-movement rule to Tnx0Vnx1Pnx2

value — then the new trees will have the same name as the input, otherwise, the name of the input tree followed by a dash ('-') and the name of the right hand side of the tree¹⁴.

The value of the parameter can be changed by choosing *Tools* at the menu bar and then

¹⁴the reason why we do not use the name of the metarule, i.e. the name of the buffer, is because in some forms of application the metarules do not carry individual names, as we'll see soon is the case when a set of metarules from a file is applied.

either *name mr output trees = input* or *append rhs name to mr output trees*.

- For applying it to all the trees of a buffer, click in the name of the buffer that contains the trees and proceed as above. The output will be a new buffer with all the output trees. The name of the new buffer will be the same as the input buffer prefixed by "MR-". The names of the trees follow the conventions above.

The other options concern application to files (instead of buffers). Lets first define the concepts of parallel, sequential and cumulative application of metarules. One metarule file can contain more than one metarule. The first two trees, i.e., the first tree pair, form one metarule - lets call it mr_0 . Subsequent pairs in the sequence of trees define additional metarules — mr_1 , mr_2 , ..., mr_n .

- We say that a metarule file is applied in parallel to a tree (see Figure B.5) if each of the metarules is applied independently to the input generating its particular output trees¹⁵. We generalize the concept to the application in parallel of a metarule file to a tree file (with possibly more than one tree), generating all the trees as if each metarule in the metarule file was applied to each tree in the input file.

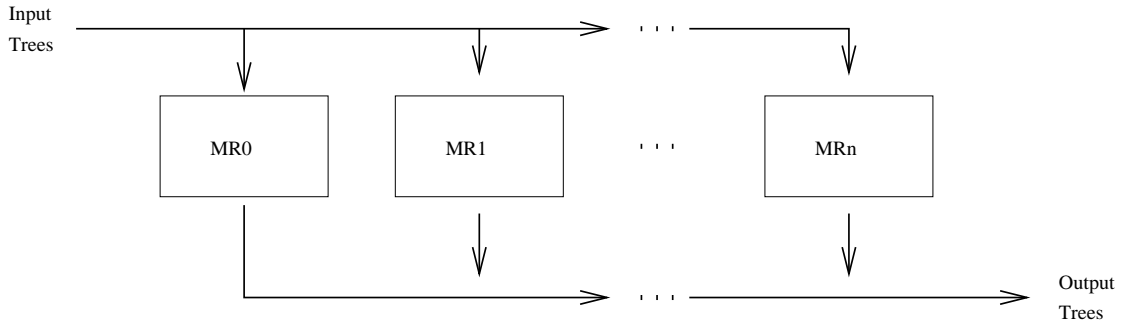


Figure B.5: Parallel application of metarules

- We say that a metarule file $mr_0, mr_1, mr_2, \dots, mr_n$ is applied in sequence to a input tree file (see Figure B.6) if we apply mr_0 to the trees of the input file, and for each $0 < i \leq n$ apply metarule mr_i to the trees generated as a result of the application of mr_{i-1} .

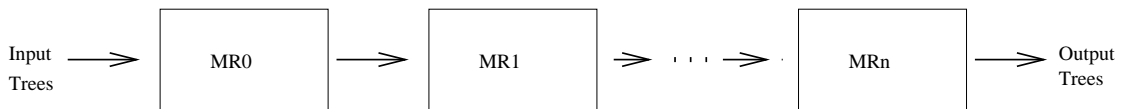


Figure B.6: Sequential application of metarules

- Finally, the cumulative application is similar to the sequential, except that the input trees at each stage are by-passed to the output together with the newly generated ones (see Figure B.7).

¹⁵remember a metarule application generates as many output trees as the number of matches

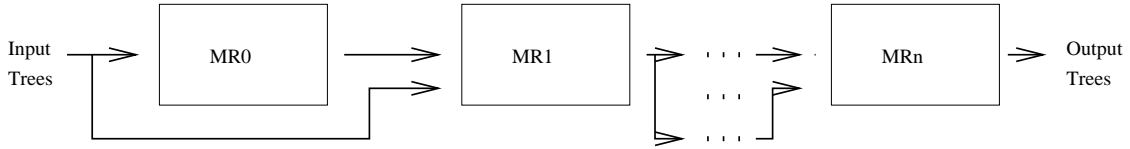


Figure B.7: Cumulative application of metarules

Remember that in case of matching failure the output result is decided as explained in subsection B.2.3 either to be empty or to be the input tree. The reflex here of having the parameter set for copying the input is that for the parallel application the output will have as many copies of the input as matching failures. For the sequential case the decision apply at each level, and setting the parameter for copying, in a certain sense, guarantees for the 'pipe' not to break. Due to its nature and unlike the two other modes, the cumulative application is not affected by this parameter.

The options for application of metarules to files are available by clicking at the menu item *Tools* and then choosing the appropriate function among:

- *Apply metarule to files:* You'll be prompted for the metarule file name which should contain one metarule¹⁶, and for input file names. Each input file name **infile** will be independently submitted to the application of the metarule generating an output file with the name **MR-infile**.
- *Apply metarules in parallel to files:* You'll be prompted for the metarules file name with one or more metarules and for input file names. Each input file name **infile** will be independently submitted to the application of the metarules in parallel. For each parallel application to a file **infile** an output file with the name **MRP-infile** will be generated.
- *Apply metarules in sequence to files:* The interaction is as described for the application in parallel, except that the application of the metarules are in sequence and that the output files are prefixed by **MRS-** instead of **MRP-**.
- *Apply metarules cumulatively to files:* The interaction is as described for the applications in parallel and in sequence, except that the mode of application is cumulative and that the output files are prefixed by **MRC-**.

Finally still under the *Tools* menu we can change the setting of the parameter that controls the output result on matching failure (see Subsection B.2.3) by choosing either *copy input on mr matching failure* or *no output on mr matching failure*.

¹⁶if it contains more than 2 trees, the additional trees are ignored

Appendix C

Lexical Organization

C.1 Introduction

An important characteristic of an FB-LTAG is that it is lexicalized, i.e., each lexical item is anchored to a tree structure that encodes subcategorization information. Trees with the same canonical subcategorizations are grouped into tree families. The reuse of tree substructures, such as *wh*-movement, in many different trees creates redundancy, which poses a problem for grammar development and maintenance [Vijay-Shanker and Schabes, 1992]. To consistently implement a change in some general aspect of the design of the grammar, all the relevant trees currently must be inspected and edited. Vijay Shanker and Schabes suggested the use of hierarchical organization and of tree descriptions to specify substructures that would be present in several elementary trees of a grammar. Since then, in addition to ourselves, Becker, [Becker, 1994], Evans et al. [Evans *et al.*, 1995], and Candito [Candito, 1996] have developed systems for organizing trees of a TAG which could be used for developing and maintaining grammars.

Our system is based on the ideas expressed in Vijay-Shanker and Schabes, [Vijay-Shanker and Schabes, 1992], to use partial-tree descriptions in specifying a grammar by separately defining pieces of tree structures to encode independent syntactic principles. Various individual specifications are then combined to form the elementary trees of the grammar. The chapter begins with a description of our grammar development system, and its implementation. We will then show the main results of using this tool to generate the Penn English grammar as well as a Chinese TAG. We describe the significant properties of both grammars, pointing out the major differences between them, and the methods by which our system is informed about these language-specific properties. The chapter ends with the conclusion and future work.

C.2 System Overview

In our approach, three types of components – subcategorization frames, blocks and lexical redistribution rules – are used to describe lexical and syntactic information. Actual trees are generated automatically from these abstract descriptions, as shown in Figure C.1. In maintaining the grammar only the abstract descriptions need ever be manipulated; the tree descriptions and the actual trees which they subsume are computed deterministically from these high-level descriptions.

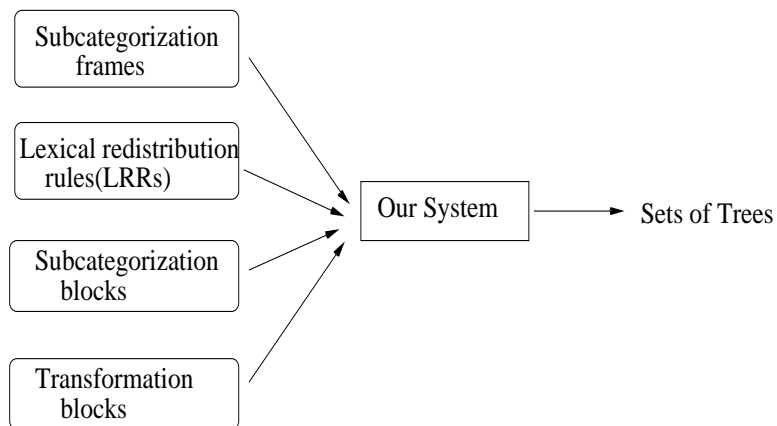


Figure C.1: Lexical Organization: System Overview

C.2.1 Subcategorization frames

Subcategorization frames specify the category of the main anchor, the number of arguments, each argument's category and position with respect to the anchor, and other information such as feature equations or node expansions. Each tree family has one canonical subcategorization frame.

C.2.2 Blocks

Blocks are used to represent the tree substructures that are reused in different trees, i.e. blocks subsume classes of trees. Each block includes a set of nodes, dominance relation, parent relation, precedence relation between nodes, and feature equations. This follows the definition of the tree descriptions specified in a logical language patterned after Rogers and Vijay-Shanker[Rogers and Vijay-Shankar, 1994].

Blocks are divided into two types according to their functions: subcategorization blocks and transformation blocks. The former describes structural configurations incorporating the various information in a subcategorization frame. For example, some of the subcategorization blocks used in the development of the English grammar are shown in Figure C.2.¹

When the subcategorization frame for a verb is given by the grammar developer, the system will automatically create a new block (of code) by essentially selecting the appropriate primitive subcategorization blocks corresponding to the argument information specified in that verb frame.

The transformation blocks are used for various transformations such as wh-movement. These transformation blocks do not encode rules for modifying trees, but rather describe the properties of a particular syntactic construction. Figure C.3 depicts our representation of phrasal extraction. This can be specialized to give the blocks for wh-movement, topicalization, relative clause formation, etc. For example, the wh-movement block is defined by further

¹In order to focus on the use of tree descriptions and to make the figures less cumbersome, we show only the structural aspects and do not show the feature value specification. The parent, (immediate dominance), relationship is illustrated by a plain line and the dominance relationship by a dotted line. The arc between nodes shows the precedence order of the nodes are unspecified. The nodes' categories are enclosed in parentheses.

APPENDIX C. LEXICAL ORGANIZATION

specifying that the ExtractionRoot is labeled S, the NewSite has a +wh feature, and so on.

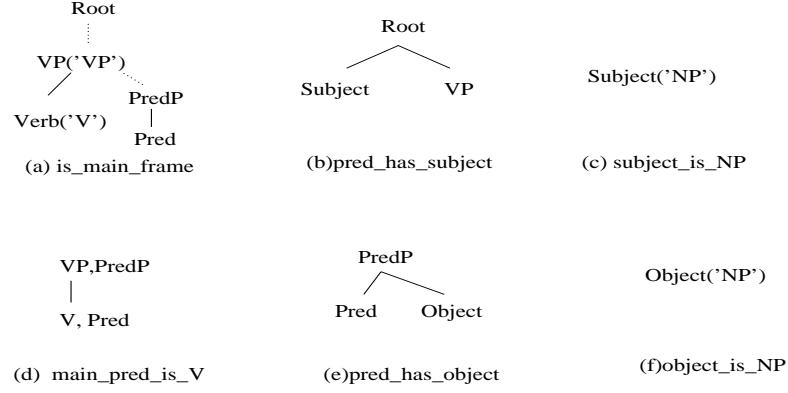


Figure C.2: Some subcategorization blocks

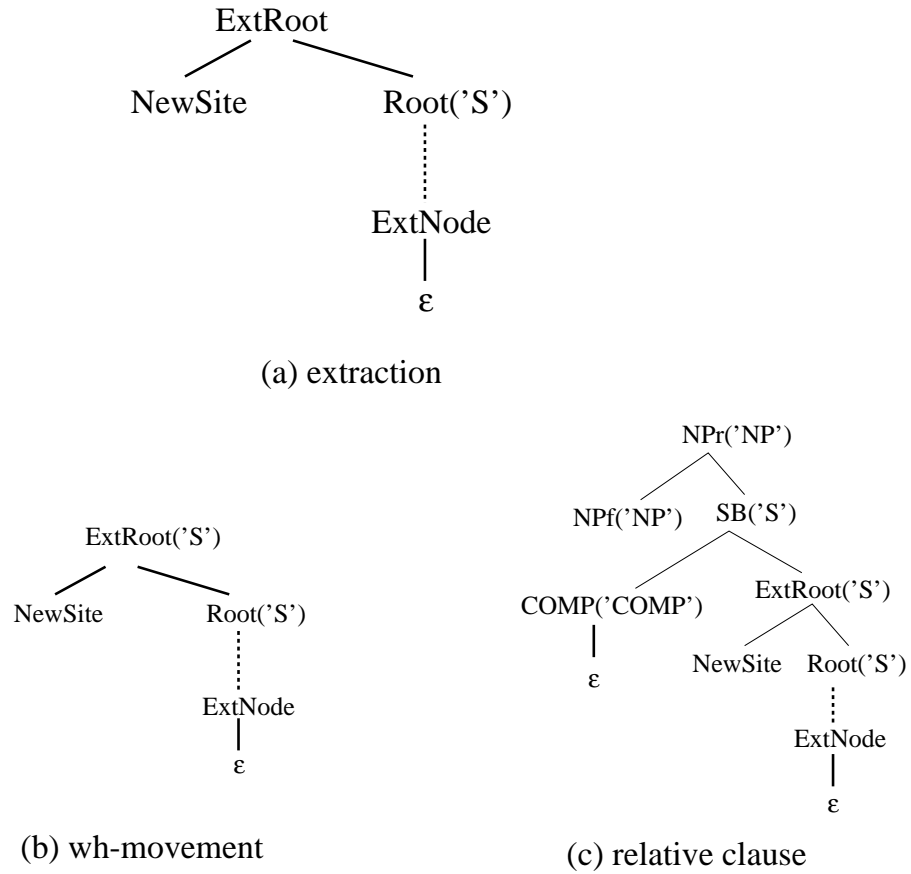


Figure C.3: Transformation blocks for extraction

C.2.3 Lexical Redistribution Rules (LRRs)

The third type of machinery available for a grammar developer is the Lexical Redistribution Rule (LRR). An LRR is a pair (r_l, r_r) of subcategorization frames, which produces a new frame when applied to a subcategorization frame s , by first *matching*² the left frame r_l of r to s , then combining information in r_r and s . LRRs are introduced to incorporate the connection between subcategorization frames. For example, most transitive verbs have a frame for active (a subject and an object) and another frame for passive, where the object in the former frame becomes the subject in the latter. An LRR, denoted as passive LRR, is built to produce the passive subcategorization frame from the active one. Similarly, applying dative-shift LRR to the frame with one NP subject and two NP objects will produce a frame with an NP subject and an PP object.

Besides the distinct content, LRRs and blocks also differ in several aspects:

- They have different functionalities: Blocks represent the substructures that are reused in different trees. They are used to reduce the redundancy among trees; LRRs are introduced to incorporate the connections between the closely related subcategorization frames.
- Blocks are strictly additive and can be added in any order. LRRs, on the other hand, produce different results depending on the order they are applied in, and are allowed to be non-additive, i.e., to remove information from the subcategorization frame they are being applied to, as in the procedure of passive from active.

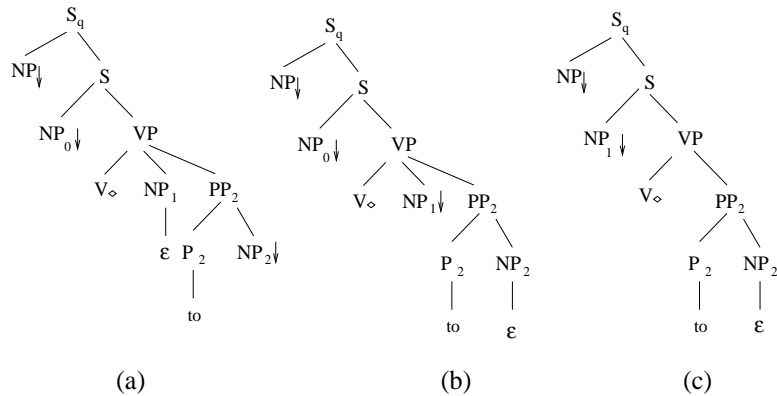


Figure C.4: Elementary trees generated from combining blocks

C.2.4 Tree generation

To generate elementary trees, we begin with a canonical subcategorization frame. The system will first generate related subcategorization frames by applying LRRs, then select subcategorization blocks corresponding to the information in the subcategorization frames, next the combinations of these blocks are further combined with the blocks corresponding to various

²Matching occurs successfully when frame s is compatible with r_l in the type of anchors, the number of arguments, their positions, categories and features. In other words, incompatible features etc. will block certain LRRs from being applied.

transformations, finally, a set of trees are generated from those combined blocks, and they are the tree family for this subcategorization frame. Figure C.4 shows some of the trees produced in this way. For instance, the last tree is obtained by incorporating information from the ditransitive verb subcategorization frame, applying the dative-shift and passive LRRs, and then combining them with the wh-non-subject extraction block. Besides, in our system the hierarchy for subcategorization frames is implicit as shown in Figure C.5.

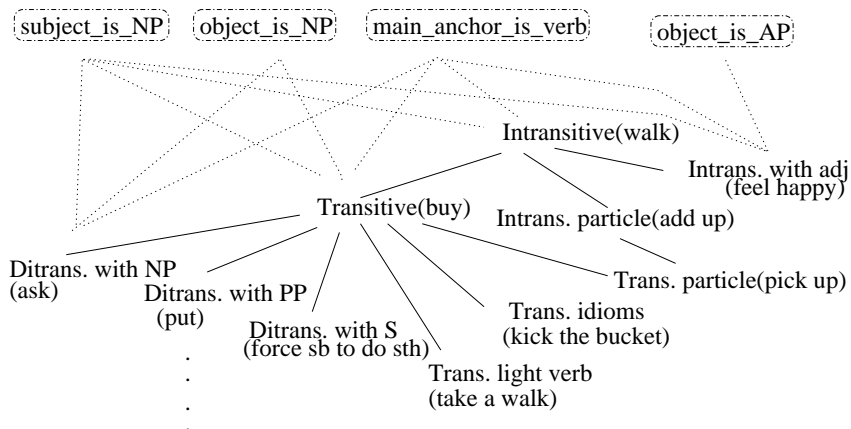


Figure C.5: Partial inheritance lattice in English

C.3 Implementation

The input of our system is the description of the language, which includes the subcategorization frame list, LRR list, subcategorization block list and transformation lists. The output is a list of trees generated automatically by the system, as shown in Figure C.6. The tree generation module is written in Prolog, and the rest part is in C. We also have a graphic interface to input the language description. Figure C.7 and C.8 are two snapshots of the interface.

C.4 Generating grammars

We have used our tool to specify a grammar for English in order to produce the trees used in the current English XTAG grammar. We have also used our tool to generate a large grammar for Chinese. In designing these grammars, we have tried to specify the grammars to reflect the similarities and the differences between the languages. The major features of our specification of these two grammars³ are summarized in Table C.1 and C.2.

³Both grammars are still under development, so the contents of these two tables might change a lot in the future according to the analyses we choose for certain phenomenon. For example, the majority of work on Chinese grammar treat *ba*-construction as some kind of object-fronting where the character *ba* is either an object marker or a preposition. According to this analysis, an LRR rule for *ba*-construction is used in our grammar to generate the preverbal-object frame from the postverbal frame. However, there has been some argument for treating *ba* as a verb. If we later choose that analysis, the main verbs in the patterns “NP0 VP” and “NP0 *ba* NP1 VP” will be different, therefore no LRR will be needed for it. As a result, the numbers of LRRs, subcat frames and tree generated will change accordingly.

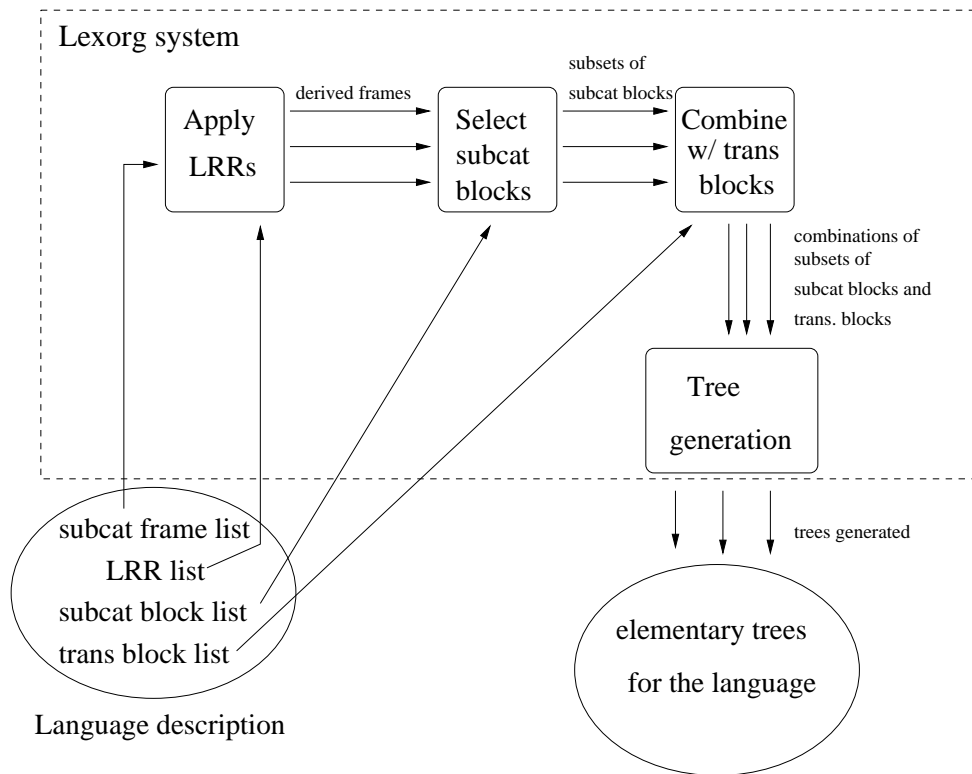


Figure C.6: Implementation of the system

File	
Block List	Block Name
Frame List	tree_has_NP_N
LRR List	is_main_frame
	main_pred_is_V
	is_main_frame
	subject_is_S
	subject_is_NP
	pred_has_subject
	pred_has_no_object
	pred_has_1_object
	pred_has_2_object
	pred_has_3_object
	object_is_S
	object1_is_S
	object2_is_S
	object_is_NP
	object1_is_NP
	object2_is_NP
	declarative
	imperative
	gerund
	tree_has_extraction_nodes
	wh_clause
	wh_subject

Figure C.7: Interface for creating a grammar

By focusing on the specification of individual grammatical information, we have been able to generate nearly all of the trees from the tree families used in the current English grammar

Block Name:

Flag: Subcategorization

Parameter List

Node List

Dom Rel List

Precedence List

Parent List

Move List

Strict Dom List

Feature Equation

Ancestor List

Node Name

Root

VP

AnchorP

Anchor

Verb

Name: VP

Type: VP

Subscript in XTAG:

☐ Parent of Lexitem
☐ Foot
☐ Subst

☐ Anchor
☐ Parent of Trace
☐ NA

☐ Extractable
☒ UnExtractable
☐ Unspecified

☐ Landable
☒ NonLandable
☐ Unspecified

Top Features:

Feature Name	Feature Value

Bottom Features:

Feature Name	Feature Value

Figure C.8: Part of the Interface for creating blocks

	English	Chinese
examples of LRRs	passive dative-shift ergative	bei-construction object fronting ba-construction
examples of transformation blocks	wh-question relativization declarative	topicalization relativization argument-drop
# LRRs	6	12
# subcat blocks	34	24
# trans blocks	8	15
# subcat frames	43	23
# trees generated	638	280

Table C.1: Major features of English and Chinese grammars

developed at Penn⁴. Our approach, has also exposed certain gaps in the Penn grammar. We are encouraged with the utility of our tool and the ease with which this large-scale grammar was developed.

We are currently working on expanding the contents of subcategorization frame to include trees for other categories of words. For example, a frame which has no specifier and one NP complement and whose predicate is a preposition will correspond to PP → P NP tree. We'll also introduce a modifier field and semantic features, so that the head features will propagate

⁴We have not yet attempted to extend our coverage to include punctuation, it-clefts, and a few idiosyncratic analyses.

	both grammars	English	Chinese
LRRs	causative short passive	long passive ergative dative-shift	VO-inversion ba-const
trans blocks	topicalization relativization declarative	gerund	argument-drop
subcat blocks	NP/S subject S/NP/PP object V predicate	PL object prep predicate	zero-subject preverbal object

Table C.2: Comparison of the two grammars

from modifiee to modified node, while non-head features from the predicate as the head of the modifier will be passed to the modified node.

C.5 Summary

We have described a tool for grammar development in which tree descriptions are used to provide an abstract specification of the linguistic phenomena relevant to a particular language. In grammar development and maintenance, only the abstract specifications need to be edited, and any changes or corrections will automatically be proliferated throughout the grammar. In addition to lightening the more tedious aspects of grammar maintenance, this approach also allows a unique perspective on the general characteristics of a language. Defining hierarchical blocks for the grammar both necessitates and facilitates an examination of the linguistic assumptions that have been made with regard to feature specification and tree-family definition. This can be very useful for gaining an overview of the theory that is being implemented and exposing gaps that remain unmotivated and need to be investigated. The type of gaps that can be exposed could include a missing subcategorization frame that might arise from the automatic combination of blocks and which would correspond to an entire tree family, a missing tree which would represent a particular type of transformation for a subcategorization frame, or inconsistent feature equations. By focusing on syntactic properties at a higher level, our approach allows new opportunities for the investigation of how languages relate to themselves and to each other.

Appendix D

Tree Naming conventions

The various trees within the XTAG grammar are named more or less according to the following tree naming conventions. Although these naming conventions are generally followed, there are occasional trees that do not strictly follow these conventions.

D.1 Tree Families

Tree families are named according to the basic declarative tree structure in the tree family (see section D.2), but with a T as the first character instead of an α or β .

D.2 Trees within tree families

Each tree begins with either an α (alpha) or a β (beta) symbol, indicating whether it is an initial or auxiliary tree, respectively. Following an α or a β the name may additionally contain one of:

I	imperative
E	ergative
N0,1,2	relative clause{position}
G	NP gerund
D	Determiner gerund
pW0,1,2	wh-PP extraction{position}
W0,1,2	wh-NP extraction{position}
X	ECM (eXceptional case marking)

Numbers are assigned according to the position of the argument in the declarative tree, as follows:

- 0 subject position
- 1 first argument (e.g. direct object)
- 2 second argument (e.g. indirect object)

The body of the name consists of a string of the following components, which corresponds to the leaves of the tree. The anchor(s) of the trees is(are) indicated by capitalizing the part of speech corresponding to the anchor.

s	sentence
a	adjective
arb	adverb
be	<i>be</i>
c	relative complementizer
x	phrasal category
d	determiner
v	verb
lv	light verb
conj	conjunction
comp	complementizer
it	<i>it</i>
n	noun
p	preposition
to	<i>to</i>
pl	particle
by	<i>by</i>
neg	negation

As an example, the transitive declarative tree consists of a subject NP, followed by a verb (which is the anchor), followed by the object NP. This translates into $\alpha nx0Vnx1$. If the subject NP had been extracted, then the tree would be $\alpha W0nx0Vnx1$. A passive tree with the *by* phrase in the same tree family would be $\alpha nx1Vbyn0$. Note that even though the object NP has moved to the subject position, it retains the object encoding (nx1).

D.3 Assorted Initial Trees

Trees that are not part of the tree families are generally gathered into several files for convenience. The various initial trees are located in `lex.trees`. All the trees in this file should begin with an α , indicating that they are initial trees. This is followed by the root category which follows the naming conventions in the previous section (e.g. n for noun, x for phrasal category). The root category is in all capital letters. After the root category, the node leaves are named, beginning from the left, with the anchor of the tree also being capitalized. As an example, the αNXN tree is rooted by an NP node (NX) and anchored by a noun (N).

D.4 Assorted Auxiliary Trees

The auxiliary trees are mostly located in the buffers `prepositions.trees`, `conjunctions.trees`, `determiners.trees`, `advs-ads.trees`, and `modifiers.trees`, although a couple of other files also contain auxiliary trees. The auxiliary trees follow a slightly different naming convention from the initial trees. Since the root and foot nodes must be the same for the auxiliary trees, the root nodes are not explicitly mentioned in the names of auxiliary trees. The trees are named according to the leaf nodes, starting from the left, and capitalizing the anchor node. All auxiliary trees begin with a β , of course. For example, $\beta ARBs$, indicates a tree anchored by

an adverb (ARB), that adjoins onto the left of an S node (Note that S must be the foot node, and therefore also the root node).

D.4.1 Relative Clause Trees

For relative clause trees, the following naming conventions have been adopted: if the *wh*-moved NP is overt, it is not explicitly represented. Instead the index of the site of movement (0 for subject, 1 for object, 2 for indirect object) is appended to the N. So $\beta N0nx0Vnx1$ is a subject extraction relative clause with \mathbf{NP}_w substitution and $\beta N1nx0Vnx1$ is an object extraction relative clause. If the *wh*-moved NP is covert and Comp substitutes in, the Comp node is represented by *c* in the tree name and the index of the extraction site follows *c*. Thus $\beta Nc0nx0Vnx1$ is a subject extraction relative clause with Comp substitution. Adjunct trees are similar, except that since the extracted material is not co-indexed to a trace, no index is specified (cf. $\beta Npxnx0Vnx1$, which is an adjunct relative clause with PP pied-piping, and $\beta Ncnx0Vnx1$, which is an adjunct relative clause with Comp substitution). Cases of pied-piping, in which the pied-piped material is part of the anchor have the anchor capitalized or spelled-out (cf. $\beta Nbynx0nx1Vbynx0$ which is a relative clause with *by*-phrase pied-piping and \mathbf{NP}_w substitution.).

Appendix E

Features

Table E.1 contains a comprehensive list of the features in the XTAG grammar and their possible values.

This section consists of short ‘biographical’ sketches of the various features currently in use in the XTAG English grammar.

E.1 Agreement

$\langle \mathbf{agr} \rangle$ is a complex feature. It can have as its subfeatures:

$\langle \mathbf{agr} \ 3\mathbf{rdsing} \rangle$, possible values: $+/-$

$\langle \mathbf{agr} \ \mathbf{num} \rangle$, possible values: *plur, sing*

$\langle \mathbf{agr} \ \mathbf{pers} \rangle$, possible values: 1, 2, 3

$\langle \mathbf{agr} \ \mathbf{gen} \rangle$, possible values: *masc, fem, neut*

These features are used to ensure agreement between a verb and its subject.

Where does it occur:

Nouns comes specified from the lexicon with their $\langle \mathbf{agr} \rangle$ features. e.g. *books* is $\langle \mathbf{agr} \ 3\mathbf{rdsing} \rangle$: $-$, $\langle \mathbf{agr} \ \mathbf{num} \rangle$: **plur**, and $\langle \mathbf{agr} \ \mathbf{pers} \rangle$: **3**. Only pronouns use the $\langle \mathbf{gen} \rangle$ (gender) feature.

The $\langle \mathbf{agr} \rangle$ features of a noun are transmitted up the NP tree by the following equation:

$$\mathbf{NP.b}:\langle \mathbf{agr} \rangle = \mathbf{N.t}:\langle \mathbf{agr} \rangle$$

Agreement between a verb and its subject is mediated by the following feature equations:

$$(509) \ \mathbf{NP}_{subj}:\langle \mathbf{agr} \rangle = \mathbf{VP.t}:\langle \mathbf{agr} \rangle$$

$$(510) \ \mathbf{VP.b}:\langle \mathbf{agr} \rangle = \mathbf{V.t}:\langle \mathbf{agr} \rangle$$

Agreement has to be done as a two step process because whether the verb agrees with the subject or not depends upon whether some auxiliary verb adjoins in and upon what the $\langle \mathbf{agr} \rangle$ specification of the verb is.

Verbs also come specified from the lexicon with their $\langle \mathbf{agr} \rangle$ features, e.g. the $\langle \mathbf{agr} \rangle$ features of the verb *sings* are $\langle \mathbf{agr} \ 3\mathbf{rdsing} \rangle$: $+$, $\langle \mathbf{agr} \ \mathbf{num} \rangle$: **sing**, and $\langle \mathbf{agr} \ \mathbf{pers} \rangle$: **3**; Non-finite forms of the verb *sing* e.g. *singing* do not come with an $\langle \mathbf{agr} \rangle$ feature specification.

Feature	Value
<agr 3rdsing>	+, -
<agr num>	plur,sing
<agr pers>	1,2,3
<agr gen>	fem,masc,neuter
<assign-case>	nom,acc,none
<assign-comp>	that,whether,if,for,ecm,rel,inf_nil,ind_nil,ppart_nil,none
<card>	+, -
<case>	nom,acc,gen,none
<comp>	that,whether,if,for,rel,inf_nil,ind_nil,nil
<compar>	+, -
<compl>	+, -
<conditional>	+, -
<conj>	and,or,but,comma,scolon,to,disc,nil
<const>	+, -
<contr>	+, -
<control>	no value, indexing only
<decreas>	+, -
<definite>	+, -
<displ-const>	+, -
<equiv>	+, -
<extracted>	+, -
<gen>	+, -
<gerund>	+, -
<inv>	+, -
<invlink>	no value, indexing only
<irrealis>	+, -
<mainv>	+, -
<mode>	base,ger,ind,inf,imp,nom,ppart,prep,subjunt
<neg>	+, -
<passive>	+, -
<perfect>	+, -
<pred>	+, -
<progressive>	+, -
<pron>	+, -
<punct bal>	dquote,squote,paren,nil
<punct contains colon>	+, -
<punct contains dash>	+, -
<punct contains dquote>	+, -
<punct contains scolon>	+, -
<punct contains squote>	+, -
<punct struct>	comma,dash,colon,scolon,nil
<punct term>	per,qmark,excl,nil
<quan>	+, -
<refl>	+, -
<rel-clause>	+, -
<rel-pron>	ppart,ger,adj-clause
<select-mode>	ind,inf,ppart,ger
<super>	+, -
<tense>	pres,past
<trace>	no value, indexing only
<trans>	+, -
<weak>	+, -
<wh>	+, -

Table E.1: List of features and their possible values

E.1.1 Agreement and Movement

The $\langle \mathbf{agr} \rangle$ features of a moved NP and its trace are co-indexed. This captures the fact that movement does not disrupt a pre-existing agreement relationship between an NP and a verb.

(511) [Which boys]_i does John think [t_i are/*is intelligent]?

E.2 Case

There are two features responsible for case-assignment:

$\langle \mathbf{case} \rangle$, possible values: **nom**, **acc**, **gen**, **none**

$\langle \mathbf{assign-case} \rangle$, possible values: **nom**, **acc**, **none**

Case assigners (prepositions and verbs) as well as the VP, S and PP nodes that dominate them have an $\langle \mathbf{assign-case} \rangle$ case feature. Phrases and lexical items that have case i.e. Ns and NPs have a $\langle \mathbf{case} \rangle$ feature.

Case assignment by prepositions involves the following equations:

$$(512) \mathbf{PP.b}:\langle \mathbf{assign-case} \rangle = \mathbf{P.t}:\langle \mathbf{case} \rangle$$

$$(513) \mathbf{NP.t}:\langle \mathbf{case} \rangle = \mathbf{P.t}:\langle \mathbf{case} \rangle$$

Prepositions come specified from the lexicon with their $\langle \mathbf{assign-case} \rangle$ feature.

$$(514) \mathbf{P.b}:\langle \mathbf{assign-case} \rangle = \mathbf{acc}$$

Case assignment by verbs has two parts: assignment of case to the object(s) and assignment of case to the subject. Assignment of case to the object is simpler. English verbs always assign accusative case to their NP objects (direct or indirect). Hence this is built into the tree and not put into the lexical entry of each individual verb.

$$(515) \mathbf{NP}_{object.t}:\langle \mathbf{case} \rangle = \mathbf{acc}$$

Assignment of case to the subject involves the following two equations.

$$(516) \mathbf{NP}_{subj}:\langle \mathbf{case} \rangle = \mathbf{VP.t}:\langle \mathbf{assign-case} \rangle$$

$$(517) \mathbf{VP.b}:\langle \mathbf{assign-case} \rangle = \mathbf{V.t}:\langle \mathbf{assign-case} \rangle$$

This is a two step process – the final case assigned to the subject depends upon the $\langle \mathbf{assign-case} \rangle$ feature of the verb as well as whether an auxiliary verb adjoins in.

Finite verbs like *sings* have **nom** as the value of their $\langle \mathbf{assign-case} \rangle$ feature. Non-finite verbs have **none** as the value of their $\langle \mathbf{assign-case} \rangle$ feature. So if no auxiliary adjoins in, the only subject they can have is **PRO** which is the only NP with **none** as the value its $\langle \mathbf{case} \rangle$ feature.

E.2.1 ECM

Certain verbs e.g. *want*, *believe*, *consider* etc. and one complementizer *for* are able to assign case to the subject of their complement clause.

The complementizer *for*, like the preposition *for*, has the $\langle \textbf{assign-case} \rangle$ feature of its complement set to **acc**. Since the $\langle \textbf{assign-case} \rangle$ feature of the root S_r of the complement tree and the $\langle \textbf{case} \rangle$ feature of its NP subject are co-indexed, this leads to the subject being assigned accusative case.

ECM verbs have the $\langle \textbf{assign-case} \rangle$ feature of their foot S node set to **acc**. The co-indexation between the $\langle \textbf{assign-case} \rangle$ feature of the root S_r and the $\langle \textbf{case} \rangle$ feature of the NP subject leads to the subject being assigned accusative case.

E.2.2 Agreement and Case

The $\langle \textbf{case} \rangle$ features of a moved NP and its trace are co-indexed. This captures the fact that movement does not disrupt a pre-existing relationship of case-assignment between a verb and an NP.

(518) $\text{Her}_i / * \text{She}_i$, I think that Odo like t_i .

E.3 Extraction and Inversion

$\langle \textbf{extracted} \rangle$, possible values are $+/-$

All sentential trees with extracted components, with the exception of relative clauses are marked **S.b** $\langle \textbf{extracted} \rangle = +$ at their top S node. The extracted element may be a *wh*-NP or a topicalized NP. The $\langle \textbf{extracted} \rangle$ feature is currently used to block embedded topicalizations as exemplified by the following example.

(519) * John wants [Bill_i [PRO to leave t_i]]

$\langle \textbf{trace} \rangle$: this feature is not assigned any value and is used to co-index moved NPs and their traces which are marked by ϵ .

$\langle \textbf{wh} \rangle$: possible values are $+/-$

NPs like *who*, *what* etc. come marked from the lexicon with a value of $+$ for the feature $\langle \textbf{wh} \rangle$. Non *wh*-NPs have $-$ as the value of their $\langle \textbf{wh} \rangle$ feature. Note that $\langle \textbf{wh} \rangle = +$ NPs are not restricted to occurring in extracted positions, to allow for the correct treatment of echo questions.

The $\langle \textbf{wh} \rangle$ feature is propagated up by possessives – e.g. the $+$ $\langle \textbf{wh} \rangle$ feature of the determiner *which* in *which boy* is propagated up to the level of the NP so that the value of the $\langle \textbf{wh} \rangle$ feature of the entire NP is $+\langle \textbf{wh} \rangle$. This process is recursive e.g. *which boy's mother*, *which boy's mother's sister*.

The $\langle \textbf{wh} \rangle$ feature is also propagated up PPs. Thus the PP *to whom* has $+$ as the value of its $\langle \textbf{wh} \rangle$ feature.

In trees with extracted NPs, the $\langle \textbf{wh} \rangle$ feature of the root node S node is equated with the $\langle \textbf{wh} \rangle$ feature of the extracted NPs.

The $\langle \textbf{wh} \rangle$ feature is used to impose subcategorizational constraints. Certain verbs like *wonder* can only take interrogative complements, other verbs such as *know* can take both

interrogative and non-interrogative complements, and yet other verbs like *think* can only take non-interrogative complements (cf. the $\langle \mathbf{extracted} \rangle$ and $\langle \mathbf{mode} \rangle$ features also play a role in imposing subcategorizational constraints).

The $\langle \mathbf{wh} \rangle$ feature is also used to get the correct inversion patterns.

E.3.1 Inversion, Part 1

The following three features are used to ensure the correct pattern of inversion:

$\langle \mathbf{wh} \rangle$: possible values are $+/-$

$\langle \mathbf{inv} \rangle$: possible values are $+/-$

$\langle \mathbf{invlink} \rangle$: possible values are $+/-$

Facts to be captured:

1. No inversion with topicalization
2. No inversion with matrix extracted subject *wh*-questions
3. Inversion with matrix extracted object *wh*-questions
4. Inversion with all matrix *wh*-questions involving extraction from an embedded clause
5. No inversion in embedded questions
6. No matrix subject topicalizations.

Consider a tree with object extraction, where NP is extracted. The following feature equations are used:

$$(520) \quad \mathbf{S}_q.\mathbf{b}:\langle \mathbf{wh} \rangle = \mathbf{NP.t}:\langle \mathbf{wh} \rangle$$

$$(521) \quad \mathbf{S}_q.\mathbf{b}:\langle \mathbf{invlink} \rangle = \mathbf{S}_q.\mathbf{b}:\langle \mathbf{inv} \rangle$$

$$(522) \quad \mathbf{S}_q.\mathbf{b}:\langle \mathbf{inv} \rangle = \mathbf{S}_r.\mathbf{t}:\langle \mathbf{inv} \rangle$$

$$(523) \quad \mathbf{S}_r.\mathbf{b}:\langle \mathbf{inv} \rangle = -$$

Root restriction: A restriction is imposed on the final root node of any XTAG derivation of a tensed sentence which equates the $\langle \mathbf{wh} \rangle$ feature and the $\langle \mathbf{invlink} \rangle$ feature of the final root node.

If the extracted NP is not a *wh*-word i.e. its $\langle \mathbf{wh} \rangle$ feature has the value $-$, at the end of the derivation, $\mathbf{S}_q.\mathbf{b}:\langle \mathbf{wh} \rangle$ will also have the value $-$. Because of the root constraint $\mathbf{S}_q.\mathbf{b}:\langle \mathbf{wh} \rangle$ will be equated to $\mathbf{S}_q.\mathbf{b}:\langle \mathbf{invlink} \rangle$ which will also come to have the value $-$. Then, by (522), $\mathbf{S}_r.\mathbf{t}:\langle \mathbf{inv} \rangle$ will acquire the value $-$. This will unify with $\mathbf{S}_r.\mathbf{b}:\langle \mathbf{inv} \rangle$ which has the value $-$ (cf. 523). Consequently, no auxiliary verb adjunction will be forced. Hence, there will never be inversion in topicalization.

If the extracted NP is a *wh*-word i.e. its $\langle \mathbf{wh} \rangle$ feature has the value $+$, at the end of the derivation, $\mathbf{S}_q.\mathbf{b}:\langle \mathbf{wh} \rangle$ will also have the value $+$. Because of the root constraint $\mathbf{S}_q.\mathbf{b}:\langle \mathbf{wh} \rangle$ will be equated to $\mathbf{S}_q.\mathbf{b}:\langle \mathbf{invlink} \rangle$ which will also come to have the value $+$. Then, by (522), $\mathbf{S}_r.\mathbf{t}:\langle \mathbf{inv} \rangle$ will acquire the value $+$. This will not unify with $\mathbf{S}_r.\mathbf{b}:\langle \mathbf{inv} \rangle$ which has the value $+$ (cf. 523). Consequently, the adjunction of an inverted auxiliary verb is required for the derivation to succeed.

Inversion will still take place even if the extraction is from an embedded clause.

(524) Who_i does Loida think [Miguel likes t_i]

This is because the adjoined tree's root node will also have its $S_r.b:\langle\mathbf{inv}\rangle$ set to $-$.

Note that inversion is only forced upon us because S_q is the final root node and the **Root restriction** applies. In embedded environments, the root restriction would not apply and the feature clash that forces adjunction would not take place.

The $\langle\mathbf{invlink}\rangle$ feature is not present in subject extractions. Consequently there is no inversion in subject questions.

Subject topicalizations are blocked by setting the $\langle\mathbf{wh}\rangle$ feature of the extracted NP to $+$ i.e. only *wh*-phrases can go in this location.

E.3.2 Inversion, Part 2

$\langle\mathbf{displ-const}\rangle$:

Possible values: $[\mathbf{set1}: +]$, $[\mathbf{set1}: -]$

In the previous section, we saw how inversion is triggered using the $\langle\mathbf{invlink}\rangle$, $\langle\mathbf{inv}\rangle$, $\langle\mathbf{wh}\rangle$ features. Inversion involves movement of the verb from V to C. This movement process is represented using the $\langle\mathbf{displ-const}\rangle$ feature which is used to simulate Multi-Component TAGs.¹ The sub-value **set1** indicates the inversion multi-component set; while there are not currently any other uses of this mechanism, it could be expanded with other sets receiving different **set** values.

The $\langle\mathbf{displ-const}\rangle$ feature is used to ensure adjunction of two trees, which in this case are the auxiliary tree corresponding to the moved verb (S adjunct) and the auxiliary tree corresponding to the trace of the moved verb (VP adjunct). The following equations are used:

$$(525) S_r.b:\langle\mathbf{displ-const set1}\rangle = -$$

$$(526) S.t:\langle\mathbf{displ-const set1}\rangle = +$$

$$(527) VP.b:\langle\mathbf{displ-const set1}\rangle = V.t:\langle\mathbf{displ-const set1}\rangle$$

$$(528) V.b:\langle\mathbf{displ-const set1}\rangle = +$$

$$(529) S_r.b:\langle\mathbf{displ-const set1}\rangle = VP.t:\langle\mathbf{displ-const set1}\rangle$$

E.4 Clause Type

There are several features that mark clause type.² They are:

$\langle\mathbf{mode}\rangle$

$\langle\mathbf{passive}\rangle$: possible values are $+/-$

$\langle\mathbf{mode}\rangle$: possible values are **base**, **ger**, **ind**, **inf**, **imp**, **nom**, **ppart**, **prep**, **sbjnt**

The $\langle\mathbf{mode}\rangle$ feature of a verb in its root form is **base**. The $\langle\mathbf{mode}\rangle$ feature of a verb in its past participial form is **ppart**, the $\langle\mathbf{mode}\rangle$ feature of a verb in its progressive/gerundive form

¹The $\langle\mathbf{displ-const}\rangle$ feature is also used in the ECM analysis.

²We have already seen one instance of a feature that marks clause-type: $\langle\mathbf{extracted}\rangle$, which marks whether a certain S involves extraction or not.

is **ger**, the $\langle \mathbf{mode} \rangle$ feature of a tensed verb is **ind**, and the $\langle \mathbf{mode} \rangle$ feature of a verb in the imperative is **imp**.

nom is the $\langle \mathbf{mode} \rangle$ value of AP/NP predicative trees headed by a null copula. **prep** is the $\langle \mathbf{mode} \rangle$ value of PP predicative trees headed by a null copula. Only the copula auxiliary tree, some sentential complement verbs (such as *consider* and raising verb auxiliary trees have **nom/prep** as the $\langle \mathbf{mode} \rangle$ feature specification of their foot node. This allow them, and only them, to adjoin onto AP/NP/PP predicative trees with null copulas.

E.4.1 Auxiliary Selection

The $\langle \mathbf{mode} \rangle$ feature is also used to state the subcategorizational constraints between an auxiliary verb and its complement. We model the following constraints:

have takes past participial complements

passive *be* takes past participial complements

active *be* takes progressive complements

modal verbs, *do*, and *to* take VPs headed by verbs in their base form as their complements.

An auxiliary verb transmits its own mode to its root and imposes its subcategorizational restrictions on its complement i.e. on its foot node. e.g. the auxiliary *have* in its infinitival form involves the following equations:

$$(530) \quad \mathbf{VP}_r.\mathbf{b}:\langle \mathbf{mode} \rangle = \mathbf{V.t}:\langle \mathbf{mode} \rangle$$

$$(531) \quad \mathbf{V.t}:\langle \mathbf{mode} \rangle = \mathbf{base}$$

$$(532) \quad \mathbf{VP.b}:\langle \mathbf{mode} \rangle = \mathbf{ppart}$$

$\langle \mathbf{passive} \rangle$: This feature is used to ensure that passives only have *be* as their auxiliary. Passive trees start out with their $\langle \mathbf{passive} \rangle$ feature as $+$. This feature starts out at the level of the verb and is percolated up to the level of the VP. This ensures that only auxiliary verbs whose foot node has $+$ as their $\langle \mathbf{passive} \rangle$ feature can adjoin on a passive. Passive trees have **ppart** as the value of their $\langle \mathbf{mode} \rangle$ feature. So the only auxiliary trees that we really have to worry about blocking are trees whose foot nodes have **ppart** as the value of their $\langle \mathbf{mode} \rangle$ feature. There are two such trees – the *be* tree and the *have* tree. The *be* tree is fine because its foot node has $+$ as its $\langle \mathbf{passive} \rangle$ feature, so both the $\langle \mathbf{passive} \rangle$ and $\langle \mathbf{mode} \rangle$ values unify; the *have* tree is blocked because its foot node has $-$ as its $\langle \mathbf{passive} \rangle$ feature.

E.5 Relative Clauses

Features that are peculiar to the relative clause system are:

$\langle \mathbf{select-mode} \rangle$, possible values are **ind**, **inf**, **ppart**, **ger**

$\langle \mathbf{rel-pron} \rangle$, possible values are **ppart**, **ger**, **adj-clause**

$\langle \mathbf{rel-clause} \rangle$, possible values are $+/ -$

$\langle \mathbf{select-mode} \rangle$:

Comps are lexically specified for $\langle \mathbf{select-mode} \rangle$. In addition, the $\langle \mathbf{select-mode} \rangle$ feature of a Comp is equated to the $\langle \mathbf{mode} \rangle$ feature of its sister S node by the following equation:

$$(533) \text{ Comp.t:}\langle\text{select-mode}\rangle = \text{S}_t.\text{t:}\langle\text{mode}\rangle$$

The lexical specifications of the Comps are shown below:

- ϵ_C , $\text{Comp.t:}\langle\text{select-mode}\rangle = \text{ind/inf/ger/ppart}$
- *that*, $\text{Comp.t:}\langle\text{select-mode}\rangle = \text{ind}$
- *for*, $\text{Comp.t:}\langle\text{select-mode}\rangle = \text{inf}$

$\langle\text{rel-pron}\rangle$:

There are additional constraints on where the null Comp ϵ_C can occur. The null Comp is not permitted in cases of subject extraction unless there is an intervening clause or the relative clause is a reduced relative (**mode = ppart/ger**).

To model this paradigm, the feature $\langle\text{rel-pron}\rangle$ is used in conjunction with the following equations.

$$(534) \text{ S}_r.\text{t:}\langle\text{rel-pron}\rangle = \text{Comp.t:}\langle\text{rel-pron}\rangle$$

$$(535) \text{ S}_r.\text{b:}\langle\text{rel-pron}\rangle = \text{S}_r.\text{b:}\langle\text{mode}\rangle$$

$$(536) \text{ Comp.b:}\langle\text{rel-pron}\rangle = \text{ppart/ger/adj-clause (for } \epsilon_C)$$

The full set of the equations above is only present in Comp substitution trees involving subject extraction. So the following will not be ruled out.

$$(537) \text{ the toy } [\epsilon_i [\epsilon_C [\text{Dafna likes } t_i]]]$$

The feature mismatch induced by the above equations is not remedied by adjunction of just any S-adjunct because all other S-adjuncts are transparent to the $\langle\text{rel-pron}\rangle$ feature because of the following equation:

$$(538) \text{ S}_m.\text{b:}\langle\text{rel-pron}\rangle = \text{S}_f.\text{t:}\langle\text{rel-pron}\rangle$$

$\langle\text{rel-clause}\rangle$:

The XTAG analysis forces the adjunction of the determiner below the relative clause. This is done by using the $\langle\text{rel-clause}\rangle$ feature. The relevant equations are:

$$(539) \text{ On the root of the RC: } \text{NP}_r.\text{b:}\langle\text{rel-clause}\rangle = +$$

$$(540) \text{ On the foot node of the Determiner tree: } \text{NP}_f.\text{t:}\langle\text{rel-clause}\rangle = -$$

E.6 Complementizer Selection

The following features are used to ensure the appropriate distribution of complementizers:

⟨**comp**⟩, possible values: **that**, **if**, **whether**, **for**, **rel**, **inf_nil**, **ind_nil**, **nil**

⟨**assign-comp**⟩, possible values: **that**, **if**, **whether**, **for**, **ecm**, **rel**, **ind_nil**, **inf_nil**, **none**

⟨**mode**⟩, possible values: **ind**, **inf**, **subjct**, **ger**, **base**, **ppart**, **nom**, **prep**

⟨**wh**⟩, possible values: **+**, **-**

The value of the ⟨**comp**⟩ feature tells us what complementizer we are dealing with. The trees which introduce complementizers come specified from the lexicon with their ⟨**comp**⟩ feature and ⟨**assign-comp**⟩ feature. The ⟨**comp**⟩ of the Comp tree regulates what kind of tree goes above the Comp tree, while the ⟨**assign-comp**⟩ feature regulates what kind of tree goes below. e.g. the following equations are used for *that*:

$$(541) \text{ } S_c.b:\langle \mathbf{comp} \rangle = \mathbf{Comp.t}:\langle \mathbf{comp} \rangle$$

$$(542) \text{ } S_c.b:\langle \mathbf{wh} \rangle = \mathbf{Comp.t}:\langle \mathbf{wh} \rangle$$

$$(543) \text{ } S_c.b:\langle \mathbf{mode} \rangle = \mathbf{ind}/\mathbf{subjct}$$

$$(544) \text{ } S_r.t:\langle \mathbf{assign-comp} \rangle = \mathbf{Comp.t}:\langle \mathbf{comp} \rangle$$

$$(545) \text{ } S_r.b:\langle \mathbf{comp} \rangle = \mathbf{nil}$$

By specifying $S_r.b:\langle \mathbf{comp} \rangle = \mathbf{nil}$, we ensure that complementizers do not adjoin onto other complementizers. The root node of a complementizer tree always has its ⟨**comp**⟩ feature set to a value other than **nil**.

Trees that take clausal complements specify with the ⟨**comp**⟩ feature on their foot node what kind of complementizer(s) they can take. The ⟨**assign-comp**⟩ feature of an S node is determined by the highest VP below the S node and the syntactic configuration the S node is in.

E.6.1 Verbs with object sentential complements

Finite sentential complements:

$$(546) \text{ } S_1.t:\langle \mathbf{comp} \rangle = \mathbf{that}/\mathbf{whether}/\mathbf{if}/\mathbf{nil}$$

$$(547) \text{ } S_1.t:\langle \mathbf{mode} \rangle = \mathbf{ind}/\mathbf{subjct} \text{ or } S_1.t:\langle \mathbf{mode} \rangle = \mathbf{ind}$$

$$(548) \text{ } S_1.t:\langle \mathbf{assign-comp} \rangle = \mathbf{ind_nil}/\mathbf{inf_nil}$$

The presence of an overt complementizer is optional.

Non-finite sentential complements, do not permit *for*:

$$(549) \text{ } S_1.t:\langle \mathbf{comp} \rangle = \mathbf{nil}$$

$$(550) \text{ } S_1.t:\langle \mathbf{mode} \rangle = \mathbf{inf}$$

(551) $S_1.t:\langle \text{assign-comp} \rangle = \text{ind_nil}/\text{inf_nil}$

Non-finite sentential complements, permit *for*:

(552) $S_1.t:\langle \text{comp} \rangle = \text{for}/\text{nil}$

(553) $S_1.t:\langle \text{mode} \rangle = \text{inf}$

(554) $S_1.t:\langle \text{assign-comp} \rangle = \text{ind_nil}/\text{inf_nil}$

Cases like ‘I want for to win’ are independently ruled out due to a case feature clash between the **acc** assigned by *for* and the intrinsic case feature **none** on the PRO.

Non-finite sentential complements, ECM:

(555) $S_1.t:\langle \text{comp} \rangle = \text{nil}$

(556) $S_1.t:\langle \text{mode} \rangle = \text{inf}$

(557) $S_1.t:\langle \text{assign-comp} \rangle = \text{ecm}$

E.6.2 Verbs with sentential subjects

The following contrast involving complementizers surfaces with sentential subjects:

(558) *(That) John is crazy is likely.

Indicative sentential subjects obligatorily have complementizers while infinitival sentential subjects may or may not have a complementizer. Also *if* is possible as the complementizer of an object clause but not as the complementizer of a sentential subject.

(559) $S_0.t:\langle \text{comp} \rangle = \text{that}/\text{whether}/\text{for}/\text{nil}$

(560) $S_0.t:\langle \text{mode} \rangle = \text{inf}/\text{ind}$

(561) $S_0.t:\langle \text{assign-comp} \rangle = \text{inf_nil}$

If the sentential subject is finite and a complementizer does not adjoin in, the $\langle \text{assign-comp} \rangle$ feature of the S_0 node of the embedding clause and the root node of the embedded clause will fail to unify. If a complementizer adjoins in, there will be no feature-mismatch because the root of the complementizer tree is not specified for the $\langle \text{assign-comp} \rangle$ feature.

The $\langle \text{comp} \rangle$ feature **nil** is split into two $\langle \text{assign-comp} \rangle$ features **ind_nil** and **inf_nil** to capture the fact that there are certain configurations in which it is acceptable for an infinitival clause to lack a complementizer but not acceptable for an indicative clause to lack a complementizer.

E.6.3 *That*-trace and *for*-trace effects

(562) Who_i do you think (*that) t_i ate the apple?

That trace violations are blocked by the presence of the following equation:

(563) $S_r.b:\langle \text{assign-comp} \rangle = \text{inf_nil}/\text{ind_nil}/\text{ecm}$

on the bottom of the S_r nodes of trees with extracted subjects (W0). The **ind_nil** feature specification permits the above example while the **inf_nil/ecm** feature specification allows the following examples to be derived:

(564) Who_i do you want [t_i to win the World Cup]?

(565) Who_i do you consider [t_i intelligent]?

The feature equation that ruled out the *that*-trace filter violations will also serve to rule out the *for*-trace violations above.

E.7 Determiner ordering

$\langle \text{card} \rangle$, possible values are +, –
 $\langle \text{compl} \rangle$, possible values are +, –
 $\langle \text{const} \rangle$, possible values are +, –
 $\langle \text{decreas} \rangle$, possible values are +, –
 $\langle \text{definite} \rangle$, possible values are +, –
 $\langle \text{gen} \rangle$, possible values are +, –
 $\langle \text{quan} \rangle$, possible values are +, –

For detailed discussion see Chapter 18.

E.8 Punctuation

$\langle \text{punct} \rangle$ is a complex feature. It has the following as its subfeatures:

$\langle \text{punct bal} \rangle$, possible values are **dquote**, **squote**, **paren**, **nil**
 $\langle \text{punct contains colon} \rangle$, possible values are +, –
 $\langle \text{punct contains dash} \rangle$, possible values are +, –
 $\langle \text{punct contains dquote} \rangle$, possible values are +, –
 $\langle \text{punct contains scolon} \rangle$, possible values are +, –
 $\langle \text{punct contains squote} \rangle$, possible values are +, –
 $\langle \text{punct struct} \rangle$, possible values are **comma**, **dash**, **colon**, **scolon**, **none**, **nil**
 $\langle \text{punct term} \rangle$, possible values are **per**, **qmark**, **excl**, **none**, **nil**

For detailed discussion see Chapter 23.

E.9 Conjunction

$\langle \mathbf{conj} \rangle$, possible values are **but**, **and**, **or**, **comma**, **scolon**, **to**, **disc**, **nil**

The $\langle \mathbf{conj} \rangle$ feature is specified in the lexicon for each conjunction and is passed up to the root node of the conjunction tree. If the conjunction is *and*, the root $\langle \mathbf{agr\ num} \rangle$ is $\langle \mathbf{plural} \rangle$, no matter what the number of the two conjuncts. With *or*, the the root $\langle \mathbf{agr\ num} \rangle$ is equated to the $\langle \mathbf{agr\ num} \rangle$ feature of the right conjunct.

The $\langle \mathbf{conj} \rangle = \mathbf{disc}$ feature is only used at the root of the β CONJs tree. It blocks the adjunction of one β CONJs tree on another. The following equations are used, where S_r is the substitution node and S_c is the root node:

$$(566) S_r.t:\langle \mathbf{conj} \rangle = \mathbf{disc}$$

$$(567) S_c.b:\langle \mathbf{conj} \rangle = \mathbf{and/or/but/nil}$$

E.10 Comparatives

$\langle \mathbf{compar} \rangle$, possible values are $+$, $-$

$\langle \mathbf{equiv} \rangle$, possible values are $+$, $-$

$\langle \mathbf{super} \rangle$, possible values are $+$, $-$

For detailed discussion see Chapter 22.

E.11 Control

$\langle \mathbf{control} \rangle$ has no value and is used only for indexing purposes. The root node of every clausal tree has its $\langle \mathbf{control} \rangle$ feature coindexed with the control feature of its subject. This allows adjunct control to take place. In addition, clauses that take infinitival clausal complements have the control feature of their subject/object coindexed with the control feature of their complement clause S, depending upon whether they are subject control verbs or object control verbs respectively.

E.12 Other Features

$\langle \mathbf{neg} \rangle$, possible values are $+$, $-$

Used for controlling the interaction of negation and auxiliary verbs.

$\langle \mathbf{pred} \rangle$, possible values are $+$, $-$

The $\langle \mathbf{pred} \rangle$ feature is used in the following tree families: Tnx0N1.trees and Tnx0nx1ARB.trees . In the Tnx0N1.trees family, the following equations are used:

for $\alpha W1nx0N1$:

$$(568) NP_1.t:\langle \mathbf{pred} \rangle = +$$

$$(569) NP_1.b:\langle \mathbf{pred} \rangle = +$$

$$(570) NP.t:\langle \mathbf{pred} \rangle = +$$

$$(571) \text{ N.t:}\langle\text{pred}\rangle = \text{NP.b:}\langle\text{pred}\rangle$$

This is the only tree in this tree family to use the $\langle\text{pred}\rangle$ feature.

The other tree family where the $\langle\text{pred}\rangle$ feature is used is Tnx0nx1ARB.trees . Within this family, this feature (and the following equations) are used only in the $\alpha\text{W1nx0nx1ARB}$ tree.

$$(572) \text{ AdvP}_1\text{.t:}\langle\text{pred}\rangle = +$$

$$(573) \text{ AdvP}_1\text{.b:}\langle\text{pred}\rangle = +$$

$$(574) \text{ NP.t:}\langle\text{pred}\rangle = +$$

$$(575) \text{ AdvP.b:}\langle\text{pred}\rangle = \text{NP.t:}\langle\text{pred}\rangle$$

$\langle\text{pron}\rangle$, possible values are $+$, $-$

This feature indicates whether a particular NP is a pronoun or not. Certain constructions which do not permit pronouns use this feature to block pronouns.

$\langle\text{tense}\rangle$, possible values are **pres**, **past**

It does not seem to be the case that the $\langle\text{tense}\rangle$ feature interacts with other features/syntactic processes. It comes from the lexicon with the verb and is transmitted up the tree in such a way that the root S node ends up with the tense feature of the highest verb in the tree. The equations used for this purpose are:

$$(576) \text{ S}_r\text{.b:}\langle\text{tense}\rangle = \text{VP.t:}\langle\text{tense}\rangle$$

$$(577) \text{ VP.b:}\langle\text{tense}\rangle = \text{V.t:}\langle\text{tense}\rangle$$

$\langle\text{trans}\rangle$, possible values are $+$, $-$

Many but not all English verbs can anchor both transitive and intransitive trees.

$$(578) \text{ The sun melted the ice cream.}$$

$$(579) \text{ The ice cream melted.}$$

$$(580) \text{ Elmo borrowed a book.}$$

$$(581) * \text{ A book borrowed.}$$

Transitive trees have the $\langle\text{trans}\rangle$ feature of their anchor set to $+$ and intransitive trees have the $\langle\text{trans}\rangle$ feature of their anchor set to $-$. Verbs such as *melt* which can occur in both transitive and intransitive trees come unspecified for the $\langle\text{trans}\rangle$ feature from the lexicon. Verbs which can only occur in transitive trees e.g. *borrow* have their $\langle\text{trans}\rangle$ feature specified in the lexicon as $+$ thus blocking their anchoring of an intransitive tree.

Appendix F

Evaluation and Results

In this appendix we describe various evaluations done of the XTAG grammar. Some of these evaluations were done on an earlier version of the XTAG grammar (the 1995 release), while other were done more recently. We will try to indicate in each section which version was used.

F.1 Parsing Corpora

In the XTAG project, we have used corpus analysis in two main ways: (1) to measure the performance of the English grammar on a given genre and (2) to identify gaps in the grammar.

The second type of evaluation involves performing detailed error analysis on the sentences rejected by the parser, and we have done this several times on WSJ and Brown data. Based on the results of such analysis, we prioritize upcoming grammar development efforts. The results of a recent error analysis are shown in Table F.1. The table does not show errors in parsing due to mistakes made by the POS tagger which contributed the largest number of errors: 32. At this point, we have added a treatment of punctuation to handle #1, an analysis of time NPs (#2), a large number of multi-word prepositions (part of #3), gapless relative clauses (#7), bare infinitives (#14) and have added the missing subcategorization (#3) and missing lexical entry (#12). We are in the process of extending the parser to handle VP coordination (#9) (See Section 21 on recent work to handle VP and other predicative coordination). We find that this method of error analysis is very useful in focusing grammar development in a productive direction.

To ensure that we are not losing coverage of certain phenomena as we extend the grammar, we have a benchmark set of grammatical and ungrammatical sentences from this technical report. We parse these sentences periodically to ensure that in adding new features and constructions to the grammar, we are not blocking previous analyses. There are approximately 590 example sentences in this set.

F.2 TSNLP

In addition to corpus-based evaluation, we have also run the English Grammar on the Test Suites for Natural Language Processing (TSNLP) English corpus [Lehmann *et al.*, 1996]. The corpus is intended to be a systematic collection of English grammatical phenomena, including

Rank	No of errors	Category of error
#1	11	Parentheticals and appositives
#2	8	Time NP
#3	8	Missing subcat
#4	7	Multi-word construction
#5	6	Ellipsis
#6	6	Not sentences
#7	3	Relative clause with no gap
#8	2	Funny coordination
#9	2	VP coordination
#10	2	Inverted predication
#11	2	Who knows
#12	1	Missing entry
#13	1	Comparative?
#14	1	Bare infinitive

Table F.1: Results of Corpus Based Error Analysis

complementation, agreement, modification, diathesis, modality, tense and aspect, sentence and clause types, coordination, and negation. It contains 1409 grammatical sentences and phrases and 3036 ungrammatical ones.

Error Class	%	Example
POS Tag	19.7%	She adds to/V it , He noises/N him abroad
Missing lex item	43.3%	<i>used</i> as an auxiliary V, <i>calm NP down</i>
Missing tree	21.2%	<i>should've</i> , <i>bet NP NP S</i> , <i>regard NP as Adj</i>
Feature clashes	3%	<i>My every firm</i> , <i>All money</i>
Rest	12.8%	<i>approx</i> , <i>e.g.</i>

Table F.2: Breakdown of TSNLP Errors

There were 42 examples which we judged ungrammatical, and removed from the test corpus. These were sentences with conjoined subject pronouns, where one or both were accusative, e.g. *Her and him succeed*. Overall, we parsed 61.4% of the 1367 remaining sentences and phrases. The errors were of various types, broken down in Table F.2. As with the error analysis described above, we used this information to help direct our grammar development efforts. It also highlighted the fact that our grammar is heavily slanted toward American English—our grammar did not handle *dare* or *need* as auxiliary verbs, and there were a number of very British particle constructions, e.g. *She misses him out*.

One general problem with the test-suite is that it uses a very restricted lexicon, and if there is one problematic lexical item it is likely to appear a large number of times and cause a disproportionate amount of grief. *Used to* appears 33 times and we got all 33 wrong. However, it must be noted that the XTAG grammar has analyses for syntactic phenomena that were not represented in the TSNLP test suite such as sentential subjects and subordinating clauses among others. This effort was, therefore, useful in highlighting some deficiencies in our grammar, but

did not provide the same sort of general evaluation as parsing corpus data.

F.3 Chunking and Dependencies in XTAG Derivations

We evaluated the XTAG parser for the text chunking task [Abney, 1991]. In particular, we compared NP chunks and verb group (VG) chunks¹ produced by the XTAG parser with the NP and VG chunks from the Penn Treebank [Marcus *et al.*, 1993]. The test involved 940 sentences of length 15 words or less from sections 17 to 23 of the Penn Treebank, parsed using the XTAG English grammar. The results are given in Table F.3.

	NP Chunking	VG Chunking
Recall	82.15%	74.51%
Precision	83.94%	76.43%

Table F.3: Text Chunking performance of the XTAG parser

System	Training Size	Recall	Precision
Ramshaw & Marcus	Baseline	81.9%	78.2%
Ramshaw & Marcus (without lexical information)	200,000	90.7%	90.5%
Ramshaw & Marcus (with lexical information)	200,000	92.3%	91.8%
Supertags	Baseline	74.0%	58.4%
Supertags	200,000	93.0%	91.8%
Supertags	1,000,000	93.8%	92.5%

Table F.4: Performance comparison of the transformation based noun chunker and the supertag based noun chunker

As described earlier, the results cannot be directly compared with other results in chunking such as in [Ramshaw and Marcus, 1995] since we do not train from the Treebank before testing. However, in earlier work, text chunking was done using a technique called supertagging [Srinivas, 1997b] (which uses the XTAG English grammar) which can be used to train from the Treebank. The comparative results of text chunking between supertagging and other methods of chunking is shown in Figure F.4.²

We also performed experiments to determine the accuracy of the derivation structures produced by XTAG on WSJ text, where the derivation tree produced after parsing XTAG is interpreted as a dependency parse. We took sentences that were 15 words or less from the Penn Treebank [Marcus *et al.*, 1993]. The sentences were collected from sections 17–23 of the Treebank. 9891 of these sentences were given at least one parse by the XTAG system. Since XTAG typically produces several derivations for each sentence we simply picked a single derivation

¹We treat a sequence of verbs and verbal modifiers, including auxiliaries, adverbs, modals as constituting a verb group.

²It is important to note in this comparison that the supertagger uses lexical information on a per word basis only to pick an initial set of supertags for a given word.

from the list for this evaluation. Better results might be achieved by ranking the output of the parser using the sort of approach described in [Srinivas *et al.*, 1995].

There were some striking differences in the dependencies implicit in the Treebank and those given by XTAG derivations. For instance, often a subject NP in the Treebank is linked with the first auxiliary verb in the tree, either a modal or a copular verb, whereas in the XTAG derivation, the same NP will be linked to the main verb. Also XTAG produces some dependencies within an NP, while a large number of words in NPs in the Treebank are directly dependent on the verb. To normalize for these facts, we took the output of the NP and VG chunker described above and accepted as correct any dependencies that were completely contained within a single chunk.

For example, for the sentence *Borrowed shares on the Amex rose to another record*, the XTAG and Treebank chunks are shown below.

XTAG chunks:

```
[Borrowed shares] [on the Amex] [rose]
[to another record]
```

Treebank chunks:

```
[Borrowed shares on the Amex] [rose]
[to another record]
```

Using these chunks, we can normalize for the fact that in the dependencies produced by XTAG *borrowed* is dependent on *shares* (i.e. in the same chunk) while in the Treebank *borrowed* is directly dependent on the verb *rose*. That is to say, we are looking at links between chunks, not between words. The dependencies for the sentence are given below.

XTAG dependency	Treebank dependency
Borrowed::shares	Borrowed::rose
shares::rose	shares::rose
on::shares	on::shares
the::Amex	the::Amex
Amex::on	Amex::on
rose::NIL	rose::NIL
to::rose	to::rose
another::record	another::record
record::to	record::to

After this normalization, testing simply consisted of counting how many of the dependency links produced by XTAG matched the Treebank dependency links. Due to some tokenization and subsequent alignment problems we could only test on 835 of the original 9891 parsed sentences. There were a total of 6135 dependency links extracted from the Treebank. The XTAG parses also produced 6135 dependency links for the same sentences. Of the dependencies produced by the XTAG parser, 5165 were correct giving us an accuracy of 84.2%.

F.4 Comparison with IBM

The evaluation in this section was done with the earlier 1995 release of the grammar. This section describes an experiment to measure the crossing bracket accuracy of the XTAG-parsed IBM-manual sentences. In this experiment, XTAG parses of 1100 IBM-manual sentences have been ranked using certain heuristics. The ranked parses have been compared³ against the bracketing given in the Lancaster Treebank of IBM-manual sentences⁴. Table F.5 shows the results of XTAG obtained in this experiment, which used the highest ranked parse for each system. It also shows the results of the latest IBM statistical grammar ([Jelinek *et al.*, 1994]) on the same genre of sentences. Only the highest-ranked parse of both systems was used for this evaluation. Crossing Brackets is the percentage of sentences with no pairs of brackets crossing the Treebank bracketing (i.e. ((a b) c) has a crossing bracket measure of one if compared to (a (b c))). Recall is the ratio of the number of constituents in the XTAG parse to the number of constituents in the corresponding Treebank sentence. Precision is the ratio of the number of correct constituents to the total number of constituents in the XTAG parse.

System	# of sentences	Crossing Bracket Accuracy	Recall	Precision
XTAG	1100	81.29%	82.34%	55.37%
IBM Statistical grammar	1100	86.20%	86.00%	85.00%

Table F.5: Performance of XTAG on IBM-manual sentences

As can be seen from Table F.5, the precision figure for the XTAG system is considerably lower than that for IBM. For the purposes of comparative evaluation against other systems, we had to use the same crossing-brackets metric though we believe that the crossing-brackets measure is inadequate for evaluating a grammar like XTAG. There are two reasons for the inadequacy. First, the parse generated by XTAG is much richer in its representation of the internal structure of certain phrases than those present in manually created treebanks (e.g. IBM: [_N your personal computer], XTAG: [_{NP} [_G your] [_N [_N personal] [_N computer]]). This is reflected in the number of constituents per sentence, shown in the last column of Table F.6.⁵

System	Sent. Length	# of sent	Av. # of words/sent	Av. # of Constituents/sent
XTAG	1-10	654	7.45	22.03
	1-15	978	9.13	30.56
IBM Stat. Grammar	1-10	447	7.50	4.60
	1-15	883	10.30	6.40

Table F.6: Constituents in XTAG parse and IBM parse

A second reason for considering the crossing bracket measure inadequate for evaluating

³We used the parseval program written by Phil Harison (phil@atc.boeing.com).

⁴The Treebank was obtained through Salim Roukos (roukos@watson.ibm.com) at IBM.

⁵We are aware of the fact that increasing the number of constituents also increases the recall percentage. However we believe that this a legitimate gain.

XTAG is that the primary structure in XTAG is the derivation tree from which the bracketed tree is derived. Two identical bracketings for a sentence can have completely different derivation trees (e.g. *kick the bucket* as an idiom vs. a compositional use). A more direct measure of the performance of XTAG would evaluate the derivation structure, which captures the dependencies between words.

F.5 Comparison with Alvey

The evaluation in this section was done with the earlier 1995 release of the grammar. This section compares XTAG to the Alvey Natural Language Tools (ANLT) Grammar. We parsed the set of LDOCE Noun Phrases presented in Appendix B of the technical report ([Carroll, 1993]) using XTAG. Table F.7 summarizes the results of this experiment. A total of 143 noun phrases were parsed. The NPs which did not have a correct parse in the top three derivations were considered failures for either system. The maximum and average number of derivations columns show the highest and the average number of derivations produced for the NPs that have a correct derivation in the top three. We show the performance of XTAG both with and without the tagger since the performance of the POS tagger is significantly degraded on the NPs because the NPs are usually shorter than the sentences on which it was trained. It would be interesting to see if the two systems performed similarly on a wider range of data.

System	# of NPs	# parsed	% parsed	Maximum derivations	Average derivations
ANLT Parser	143	127	88.81%	32	4.57
XTAG Parser with POS tagger	143	93	65.03%	28	3.45
XTAG Parser without POS tagger	143	120	83.91%	28	4.14

Table F.7: Comparison of XTAG and ANLT Parser

F.6 Comparison with CLARE

The evaluation in this section was done with the earlier 1995 release of the grammar. This section compares the performance of XTAG against that of the CLARE-2 system ([Alshaw et al., 1992]) on the ATIS corpus. Table F.8 shows the performance results. The percentage parsed column for both systems represents the percentage of sentences that produced any parse. It must be noted that the performance result shown for CLARE-2 is without any tuning of the grammar for the ATIS domain. The performance of CLARE-3, a later version of the CLARE system, is estimated to be 10% higher than that of the CLARE-2 system.⁶

In an attempt to compare the performance of the two systems on a wider range of sentences (from similar genres), we provide in Table F.9 the performance of CLARE-2 on LOB corpus and

⁶When CLARE-3 is tuned to the ATIS domain, performance increases to 90%. However XTAG has not been tuned to the ATIS domain.

APPENDIX F. EVALUATION AND RESULTS

System	Mean length	% parsed
CLARE-2	6.53	68.50%
XTAG	7.62	88.35%

Table F.8: Performance of CLARE-2 and XTAG on the ATIS corpus

the performance of XTAG on the WSJ corpus. The performance was measured on sentences of up to 10 words for both systems.

System	Corpus	Mean length	% parsed
CLARE-2	LOB	5.95	53.40%
XTAG	WSJ	6.00	55.58%

Table F.9: Performance of CLARE-2 and XTAG on LOB and WSJ corpus respectively

Bibliography

- [Abeillé and Schabes, 1989] Anne Abeillé and Yves Schabes. Parsing Idioms in Lexicalized TAGs. In *Proceedings of EACL '89*, pages 161–65, 1989.
- [Abeillé, 1990] Anne Abeillé. French and english determiners: Interaction of morphology, syntax, and semantics in lexicalized tree adjoining grammars. In *Tree Adjoining Grammar, First International Workshop on TAGs: Formal Theory and Applications (abstracts)*, Schloss Dagstuhl, Sarrbrücken, 1990.
- [Abney, 1987] Steven Abney. *The English Noun Phrase in its Sentential Aspects*. PhD thesis, MIT, 1987.
- [Abney, 1991] Steven Abney. Parsing by chunks. In Robert Berwick, Steven Abney, and Carol Tenny, editors, *Principle-based parsing*. Kluwer Academic Publishers, 1991.
- [Akmajian, 1970] A. Akmajian. On deriving cleft sentences from pseudo-cleft sentences. *Linguistic Inquiry*, 1:149–168, 1970.
- [Alshawhi *et al.*, 1992] Hiyaw Alshawhi, David Carter, Richard Crouch, Steve Pullman, Manny Rayner, and Arnold Smith. *CLARE – A Contextual Reasoning and Cooperative Response Framework for the Core Language Engine*. SRI International, Cambridge, England, 1992.
- [Ball, 1991] Catherine N. Ball. *The historical development of the it-cleft*. PhD thesis, University of Pennsylvania, 1991.
- [Baltin, 1989] Mark Baltin. Heads and projections. In Mark Baltin and Anthony S. Kroch, editors, *Alternative conceptions of phrase structure*, pages 1–16. University of Chicago Press, Chicago, Illinois, 1989.
- [Barwise and Cooper, 1981] John Barwise and Robin Cooper. Generalized Quantifiers and Natural Language. *Linguistics and Philosophy*, 4, 1981.
- [Becker, 1993] Tilman Becker. *HyTAG: A new Type of Tree Adjoining Grammars for Hybrid Syntactic Representation of Free Word Order Languages*. PhD thesis, Universität des Saarlandes, 1993.
- [Becker, 1994] Tilman Becker. Patterns in metarules. In *Proceedings of the 3rd TAG+ Conference*, Paris, France, 1994.
- [Bhatt, 1994] Rajesh Bhatt. Pro-control in tags. Unpublished paper, University of Pennsylvania, 1994.

- [Browning, 1987] Marguerite Browning. *Null Operator Constructions*. PhD thesis, MIT, 1987.
- [Burzio, 1986] Luigi Burzio. *Italian syntax. A Government-Binding approach*. Studies in natural language and linguistic theory. Reidel, Dordrecht, 1986.
- [Candito, 1996] Marie-Helene Candito. A Principle-Based Hierarchical Representation of LT-AGs. In *Proceedings of COLING-96*, Copenhagen, Denmark, 1996.
- [Carroll, 1993] John Carroll. *Practical Unification-based Parsing of Natural Language*. University of Cambridge, Computer Laboratory, Cambridge, England, 1993.
- [Chomsky and Lasnik, 1993] Noam Chomsky and Howard Lasnik. The minimalist program. manuscript, 1993.
- [Chomsky, 1965] Noam Chomsky. *Aspects of the Theory of Syntax*. MIT Press, Cambridge, Massachusetts, 1965.
- [Chomsky, 1970] Noam Chomsky. Remarks on Nominalization. In *Readings in English Transformational Grammar*, pages 184–221. Ginn and Company, Waltham, Massachusetts, 1970.
- [Chomsky, 1986] Noam Chomsky. *Barriers*. MIT Press, Cambridge, Massachusetts, 1986.
- [Chomsky, 1992] Noam Chomsky. A Minimalist Approach to Linguistic Theory. *MIT Working Papers in Linguistics*, Occasional Papers in Linguistics No. 1, 1992.
- [Church, 1988] Kenneth Ward Church. A Stochastic Parts Program and Noun Phrase Parser for Unrestricted Text. In *2nd Applied Natural Language Processing Conference*, Austin, Texas, 1988.
- [Cinque, 1990] G. Cinque. *Types of A-bar-Dependencies*. MIT Press, Cambridge, Massachusetts; London, England, 1990.
- [Cowie and Mackin, 1975] A. P. Cowie and R. Mackin, editors. *Oxford Dictionary of Current Idiomatic English*, volume 1. Oxford University Press, London, England, 1975.
- [Delahunty, 1984] Gerald P. Delahunty. The Analysis of English Cleft Sentences. *Linguistic Analysis*, 13(2):63–113, 1984.
- [Delin, 1989] Judy L. Delin. *Cleft Constructions in Discourse*. PhD thesis, University of Edinburgh, 1989.
- [Doran, 1998] Christine Doran. *Incorporating Punctuation into the Sentence Grammar: A Lexicalized Tree Adjoining Grammar Perspective*. PhD thesis, University of Pennsylvania, 1998.
- [Egedi and Martin, 1994] Dania Egedi and Patrick Martin. A Freely Available Syntactic Lexicon for English. In *Proceedings of the International Workshop on Sharable Natural Language Resources*, Nara, Japan, August 1994.
- [Emonds, 1970] J. Emonds. *Root and structure-preserving transformations*. PhD thesis, Massachusetts Institute of Technology, 1970.

- [Ernst, 1983] T. Ernst. More on Adverbs and Stressed Auxiliaries. *Linguistic Inquiry* 13, pages 542–48, 1983.
- [Evans *et al.*, 1995] Roger Evans, Gerald Gazdar, and David Weir. Encoding Lexicalized Tree Adjoining Grammars with a Nonmonotonic Inheritance Hierarchy. In *Proceedings of the 33rd Annual Meeting of the Association for Computational Linguistics*, Cambridge, MA, 1995.
- [Gazdar *et al.*, 1985] G. Gazdar, E. Klein, G. Pullum, and I. Sag. *Generalized Phrase Structure Grammar*. Harvard University Press, Cambridge, Massachusetts, 1985.
- [Grimshaw, 1990] Jane Grimshaw. *Argument Structure*. MIT Press, Cambridge, Massachusetts, 1990.
- [Haegeman, 1991] Liliane Haegeman. *Introduction to Government and Binding Theory*. Blackwell Publishers, Cambridge, Massachusetts, Oxford, U.K., 1991.
- [Hale and Keyser, 1986] Ken Hale and Samuel Jay Keyser. Some transitivity alternations in English. Technical Report #7, MITCCS, July 1986.
- [Hale and Keyser, 1987] Ken Hale and Samuel Jay Keyser. A view from the middle. Technical Report #10, MITCCS, January 1987.
- [Hanks, 1979] Patrick Hanks, editor. *Collins Dictionary of the English Language*. Collins, London, England, 1979.
- [Heycock and Kroch, 1993] Caroline Heycock and Anthony S. Kroch. Verb movement and the status of subjects: implications for the theory of licensing. *GAGL*, 36:75–102, 1993.
- [Heycock, 1991] Caroline Heycock. Progress Report: The English Copula in a LTAG. Internal XTAG project report, University of Pennsylvania, Summer 1991.
- [Hockey and Mateyak, 1998] Beth Ann Hockey and Heather Mateyak. Determining determiner sequencing: A syntactic analysis for english. Technical Report 98-10, Institute for Research in Cognitive Science, University of Pennsylvania, 1998.
- [Hockey and Srinivas, 1993] Beth Ann Hockey and B. Srinivas. Feature-Based TAG in Place of Multi-component Adjunction: Computational Implications. In *Proceedings of the Natural Language Processing Pacific Rim Symposium (NLPRS)*, Fukuoka, Japan, December 1993.
- [Hockey, 1994] Beth Ann Hockey. Echo questions, intonation and focus. In *Proceedings of the Interdisciplinary Conference on Focus and Natural Language Processing in Celebration of the 10th anniversary of the Journal of Semantics*, Eschwege, Germany, 1994.
- [Hornby, 1974] A. S. Hornby, editor. *Oxford Advanced Learner’s Dictionary of Current English*. Oxford University Press, London, England, third edition, 1974.
- [Iatridou, 1991] Sabine Iatridou. *Topics in Conditionals*. PhD thesis, MIT, 1991.
- [Jackendoff, 1972] R. Jackendoff. *Semantic Interpretation in Generative Grammar*. MIT Press, Cambridge, Massachusetts, 1972.

- [Jackendoff, 1990] R. Jackendoff. On Larson's Analysis of the Double Object Construction. *Linguistic Inquiry*, 21:427–456, 1990.
- [Jelinek *et al.*, 1994] Fred Jelinek, John Lafferty, David M. Magerman, Robert Mercer, Adwait Ratnaparkhi, and Salim Roukos. Decision Tree Parsing using a Hidden Derivation Model. In *Proceedings from the ARPA Workshop on Human Language Technology Workshop*, March 1994.
- [Joshi *et al.*, 1975] Aravind K. Joshi, L. Levy, and M. Takahashi. Tree Adjunct Grammars. *Journal of Computer and System Sciences*, 1975.
- [Joshi, 1985] Aravind K. Joshi. Tree Adjoining Grammars: How much context Sensitivity is required to provide a reasonable structural description. In D. Dowty, I. Karttunen, and A. Zwicky, editors, *Natural Language Parsing*, pages 206–250. Cambridge University Press, Cambridge, U.K., 1985.
- [Kaplan and Bresnan, 1983] Ronald Kaplan and Joan Bresnan. Lexical-functional Grammar: A Formal System for Grammatical Representation. In J. Bresnan, editor, *The Mental Representation of Grammatical Relations*. MIT Press, Cambridge, Massachusetts, 1983.
- [Karp *et al.*, 1992] Daniel Karp, Yves Schabes, Martin Zaidel, and Dania Egedi. A Freely Available Wide Coverage Morphological Analyzer for English. In *Proceedings of the 15th International Conference on Computational Linguistics (COLING '92)*, Nantes, France, August 1992.
- [Keenan and Stavi, 1986] E. L. Keenan and J. Stavi. A Semantic Characterization of Natural Language Determiners. *Linguistics and Philosophy*, 9, August 1986.
- [Knowles, 1986] John Knowles. The Cleft Sentence: A Base-Generated Perspective. *Lingua: International Review of General Linguistics*, 69(4):295–317, August 1986.
- [Kroch and Joshi, 1985] Anthony S. Kroch and Aravind K. Joshi. The Linguistic Relevance of Tree Adjoining Grammars. Technical Report MS-CIS-85-16, Department of Computer and Information Science, University of Pennsylvania, 1985.
- [Kroch and Joshi, 1987] Anthony S. Kroch and Aravind K. Joshi. Analyzing Extraposition in a Tree Adjoining Grammar. In G. Huck and A. Ojeda, editors, *Discontinuous Constituents, Syntax and Semantics*, volume 20. Academic Press, 1987.
- [Lapointe, 1980] S. Lapointe. A lexical analysis of the English auxiliary verb system. In T. Hoekstra *et al.*, editor, *Lexical Grammar*, pages 215–254. Foris, Dordrecht, 1980.
- [Larson, 1988] R. Larson. On the Double Object construction. *Linguistic Inquiry*, 19:335–391, 1988.
- [Larson, 1990] R. Larson. Double Objects Revisited: Reply to Jackendoff. *Linguistic Inquiry*, 21:589–612, 1990.
- [Lasnik and Saito, 1984] H. Lasnik and M. Saito. On the Nature of Proper Government. *Linguistic Inquiry*, 15:235–289, 1984.

- [Lasnik and Uriagereka, 1988] H. Lasnik and J. Uriagereka. *A Course in GB Syntax*. MIT Press, Cambridge, Massachusetts, 1988.
- [Lees, 1960] Robert B. Lees. *The Grammar of English Nominalizations*. Indiana University Research Center in Anthropology, Folklore, and Linguistics, Indiana University, Bloomington, Indiana, 1960.
- [Lehmann *et al.*, 1996] Sabine Lehmann, Stephan Oepen, Sylvie Regnier-Prost, Klaus Netter, Veronika Lux, Judith Klein, Kirsten Falkedal, Frederik Fouvry, Dominique Estival, Eva Dauphin, Hervé Compagnion, Judith Baur, Lorna Balkan, and Doug Arnold. TSNLP — Test Suites for Natural Language Processing. In *Proceedings of COLING 1996*, Copenhagen, 1996.
- [Lieberman, 1989] Mark Liberman. Tex on tap: the ACL Data Collection Initiative. In *Proceedings of the DARPA Workshop on Speech and Natural Language Processing*. Morgan Kaufman, 1989.
- [Marcus *et al.*, 1993] Mitchell M. Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. Building a Large Annotated Corpus of English: The Penn Treebank. *Computational Linguistics*, 19.2:313–330, June 1993.
- [Mateyak, 1997] Heather Mateyak. Negation of noun phrases with *not*. Technical Report 97-18, UPENN-IRCS, October 1997.
- [McCawley, 1988] James D. McCawley. *The Syntactic Phenomena of English*. The University of Chicago Press, Chicago, Illinois, 1988.
- [Moro, 1990] A. Moro. *There* as a Raised Predicate. Paper presented at GLOW, 1990.
- [Napoli, 1988] Donna Jo Napoli. Ergative Verbs in English. *Journal of the Linguistic Society of America (Language)*, 64:1(1):130–142, March 1988.
- [Oberbauer, 1984] H. Oberbauer. On the Identification of Empty Categories. *Linguistic Review* 4, 2:153–202, 1984.
- [Paroubek *et al.*, 1992] Patrick Paroubek, Yves Schabes, and Aravind K. Joshi. Xtag – a graphical workbench for developing tree-adjoining grammars. In *Third Conference on Applied Natural Language Processing*, Trento, Italy, 1992.
- [Partee *et al.*, 1990] Barbara Partee, Alice ter Meulen, and Robert E. Wall. *Mathematical Methods in Linguistics*. Kluwer Academic Publishers, 1990.
- [Pollard and Sag, 1994] Carl Pollard and Ivan A. Sag. *Head-Driven Phrase Structure Grammar*. CSLI, 1994.
- [Pollock, 1989] J-Y. Pollock. Verb Movement, UG, and the Structure of IP. *Linguistic Inquiry*, 20.3:365–424, 1989.
- [Quirk *et al.*, 1985] Randolph Quirk, Sidney Greenbaum, Geoffrey Leech, and Jan Svartvik. *A Comprehensive Grammar of the English Language*. Longman, London, 1985.

- [Ramshaw and Marcus, 1995] Lance Ramshaw and Mitchell P. Marcus. Text chunking using transformation-based learning. In *Proceedings of the Third Workshop on Very Large Corpora*, MIT, Cambridge, Boston, 1995.
- [Rizzi, 1990] Luigi Rizzi. *Relativized Minimality*. MIT Press, Cambridge, Massachusetts; London, England, 1990.
- [Rogers and Vijay-Shankar, 1994] James Rogers and K. Vijay-Shankar. Obtaining Trees from their Descriptions: An Application to Tree Adjoining Grammars. *Computational Intelligence*, 10(4), 1994.
- [Rosenbaum, 1967] Peter S. Rosenbaum. *The grammar of English predicate complement constructions*. MIT press, Cambridge, Massachusetts, 1967.
- [Sag *et al.*, 1985] I. Sag, G. Gazdar, T. Wasow, and S. Weisler. Coordination and How to distinguish categories. *Natural Language and Linguistic Theory*, 3:117–171, 1985.
- [Sarkar and Joshi, 1996] Anoop Sarkar and Aravind Joshi. Coordination in Tree Adjoining Grammars: Formalization and Implementation. In *Proceedings of the 18th International Conference on Computational Linguistics (COLING '94)*, Copenhagen, Denmark, August 1996.
- [Schabes and Joshi, 1988] Yves Schabes and Aravind K. Joshi. An Early-Type Parsing Algorithm for Tree Adjoining Grammars. In *Proceedings of the 26th Meeting of the Association for Computational Linguistics*, Buffalo, June 1988.
- [Schabes *et al.*, 1988] Yves Schabes, Anne Abeillé, and Aravind K. Joshi. Parsing strategies with ‘lexicalized’ grammars: Application to Tree Adjoining Grammars. In *Proceedings of the 12th International Conference on Computational Linguistics (COLING'88)*, Budapest, Hungary, August 1988.
- [Schabes, 1990] Yves Schabes. *Mathematical and Computational Aspects of Lexicalized Grammars*. PhD thesis, Computer Science Department, University of Pennsylvania, 1990.
- [Seuss, 1971] Dr. Seuss. *The Lorax*. Random House, New York, New York, 1971.
- [Soong and Huang, 1990] Frank K. Soong and Eng-Fong Huang. Fast Tree-Trellis Search for Finding the N-Best Sentence Hypothesis in Continuous Speech Recognition. *Journal of Acoustic Society, AM.*, May 1990.
- [Sornicola, 1988] R. Sornicola. IT-clefts and WH-clefts: two awkward sentence types. *Journal of Linguistics*, 24:343–379, 1988.
- [Srinivas *et al.*, 1994] B. Srinivas, D. Egedi, C. Doran, and T. Becker. Lexicalization and grammar development. In *Proceedings of KONVENS '94*, pages 310–9, Vienna, Austria, 1994.
- [Srinivas *et al.*, 1995] B. Srinivas, Christine Doran, and Seth Kulick. Heuristics and parse ranking. In *Proceedings of the 4th Annual International Workshop on Parsing Technologies*, Prague, September 1995.

- [Srinivas, 1997a] B. Srinivas. *Complexity of Lexical Descriptions and its Relevance to Partial Parsing*. PhD thesis, University of Pennsylvania, Philadelphia, PA, August 1997.
- [Srinivas, 1997b] B. Srinivas. Performance Evaluation of Supertagging for Partial Parsing. In *Proceedings of Fifth International Workshop on Parsing Technology*, Boston, USA, September 1997.
- [Vijay-Shanker and Joshi, 1991] K. Vijay-Shanker and Aravind K. Joshi. Unification Based Tree Adjoining Grammars. In J. Wedekind, editor, *Unification-based Grammars*. MIT Press, Cambridge, Massachusetts, 1991.
- [Vijay-Shanker and Schabes, 1992] K. Vijay-Shanker and Yves Schabes. Structure sharing in lexicalized tree adjoining grammar. In *Proceedings of the 15th International Conference on Computational Linguistics (COLING '92)*, Nantes, France, August 1992.
- [Watanabe, 1993] Akira Watanabe. The Notion of Finite Clauses in AGR-Based Case Theory. *MIT Working Papers in Linguistics*, 18:281–296, 1993.
- [XTAG-Group, 1995] The XTAG-Group. A Lexicalized Tree Adjoining Grammar for English. Technical Report IRCS 95-03, University of Pennsylvania, 1995.